



# Tree-Based Morse Regions: A Topological Approach to Local Feature Detection

Yongchao Xu, Pascal Monasse, Thierry Géraud, Laurent Najman

## ► To cite this version:

Yongchao Xu, Pascal Monasse, Thierry Géraud, Laurent Najman. Tree-Based Morse Regions: A Topological Approach to Local Feature Detection. IEEE Transactions on Image Processing, 2014, 23 (12), pp.5612-5625. 10.1109/TIP.2014.2364127 . hal-01162446

**HAL Id: hal-01162446**

**<https://hal.science/hal-01162446>**

Submitted on 10 Jun 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Tree-Based Morse Regions: A Topological Approach to Local Feature Detection

Yongchao Xu, Pascal Monasse, Thierry Géraud, and Laurent Najman

**Abstract**—This paper introduces a topological approach to local invariant feature detection motivated by Morse theory. We use the critical points of the graph of the intensity image, revealing directly the topology information as initial “interest” points. Critical points are selected from what we call a tree-based shape-space. Specifically, they are selected from both the connected components of the upper level sets of the image (the Max-tree) and those of the lower level sets (the Min-tree). They correspond to specific nodes on those two trees: (1) to the leaves (extrema) and (2) to the nodes having bifurcation (saddle points). We then associate to each critical point the largest region that contains it and is topologically equivalent in its tree. We call such largest regions the Tree-Based Morse Regions (TBMR).

TBMR can be seen as a variant of MSER, which are contrasted regions. Contrarily to MSER, TBMR relies only on topological information and thus fully inherit the invariance properties of the space of shapes (e.g., invariance to affine contrast changes and covariance to continuous transformations). In particular, TBMR extracts the regions independently of the contrast, which makes it truly contrast invariant. Furthermore, it is quasi parameter-free. TBMR extraction is fast, having the same complexity as MSER. Experimentally, TBMR achieves a repeatability on par with state-of-the-art methods, but obtains a significantly higher number of features. Both the accuracy and the robustness of TBMR are demonstrated by applications to image registration and 3D reconstruction.

**Index Terms**—Min/Max tree, local features, affine region detectors, image registration, 3D reconstruction.

## I. INTRODUCTION

LOCAL invariant feature detection [1, 2, 3, 4, 5, 6, 7, 8] is an important step in a number of applications such as wide baseline matching, object and image retrieval, tracking, recognition, image registration and 3D reconstruction. The classical process to obtain the features consists in detecting a specific class of interest points, such as corners, together with an associated scale generally obtained from a scale-space. Typical examples of such key locations are the local extrema of the result of difference of Gaussians (DoG) applied in scale-space to a series of smoothed and resampled images.

Yongchao Xu is with EPITA Research and Development Laboratory, 14-16 rue Voltaire, FR-94270 Le Kremlin-Bicêtre, France, and with the Laboratoire d’Informatique Gaspard-Monge, Université Paris-Est, Équipe A3SI, ESIEE Paris, 2 bd Blaise Pascal, B.P. 99, FR-93160 Noisy-le-Grand, France (e-mail: yongchao.xu@lrde.epita.fr).

Pascal Monasse is with the Laboratoire d’Informatique Gaspard-Monge, Université Paris-Est, Équipe A3SI, Imagine group at École des Ponts Paris-Tech, France (e-mail: monasse@imagine.enpc.fr).

Thierry Géraud is with EPITA Research and Development Laboratory, 14-16 rue Voltaire, FR-94270 Le Kremlin-Bicêtre, France (e-mail: thierry.geraud@lrde.epita.fr).

Laurent Najman is with the Laboratoire d’Informatique Gaspard-Monge, Université Paris-Est, Équipe A3SI, ESIEE Paris, 2 bd Blaise Pascal, B.P. 99, FR-93160 Noisy-le-Grand, France (e-mail: l.najman@esiee.fr).



(a) Four among 76 used multi-view images.



(b) Incomplete reconstructed 3D facades using DoG.



(c) The four facades of the 3D reconstruction using TBMR.

Fig. 1. An example of 3D reconstruction using local invariant features. Top: 4 among 76 used multi-view images. Four facades of the PMVS [9] densified sparse 3D reconstruction from the SfM pipeline [10] using DoG (middle: only 62 images are calibrated, and the back facade of the building is missing) and the proposed TBMR (down: all the 76 images are calibrated, and the whole mansion is reconstructed).

Several crucial invariance properties are required for using such points in applications, such as invariance to image translation, scaling, and rotation, to illumination changes or to local geometric distortion.

In this paper, we propose a topological approach to extract the local invariant features. We first extract some initial “critical” points, based on ideas from the Morse theory [11]: minima, maxima, and saddle points. More precisely, following [12], we propose to choose critical regions in the two trees (called Min-tree and Max-tree [13]) made by the connected components of lower and upper level sets: those critical regions are the leaves and the regions resulting from a fork. For each critical region a scale is selected. Instead of using a scale-space, the scale comes from the tree-based shape space which is built from the image  $f$  using the Max-tree and Min-tree: we associate to a critical region  $R_c$  the largest region containing it and topologically equivalent in its tree. We call our method Tree-Based Morse Regions (TBMR).

As detailed in Section III-A, the tree-based shape space



is invariant to affine contrast changes and to continuous (topological) transformations such as translation, scaling or rotation. As TBMR uses only topological information, TBMR inherits from this shape space, and is thus independent of the image contrast and covariant to these same continuous transformations. As demonstrated in this paper, TBMR is also robust to local geometric distortion. Furthermore, it is essentially parameter-free: only two non-significant parameters are applied, so that we ignore regions that are either too small or too large. Besides, the maximum area parameter is actually not important: large shapes are few and do not influence much the results. And last, but not the least, efficient algorithms with a quasi-linear or a linear complexity are available to compute the TBMR [13, 14, 15].

Qualitative experiments (Section V-A) show the better distribution of TBMR compared to other state-of-the-art methods. Quantitative evaluation, based on the image coverage measurement in Section V-B, confirms the qualitative evaluation. Tests in Section V-C demonstrate that TBMR achieves repeatability score comparable to other state-of-the-art methods with a significantly higher number of correspondences. We evaluate TBMR on two applications in which many matched features are required: image registration (Section V-D) and 3D reconstruction (Section V-E). For these two applications, and as illustrated in Fig. 1, results attest that TBMR improves over the commonly used DoG [16].

## II. RELATED WORK

There exist a variety of local invariant feature detectors having relatively good performance, as assessed by several evaluation frameworks [5, 6, 7, 8]. The first type is based on scale-space. Harris corners, Hessian based detectors, and the Difference of Gaussians (DoG) are such instances. The Harris corner detector [17] finds the extrema of a corner measure based on the second moment matrix at some fixed scale. A scale-adapted Harris corner detector and its extension Harris-Laplace [18] with scale selection find extrema of the Laplacian of Gaussian (LoG) filter. The Hessian detector [3] extracts the extrema of a feature measure based on the Hessian matrix. Its extension Hessian-Laplace [18] uses the same scale selection as Harris-Laplace. The affine versions of both Harris and Hessian are based on the affine shape estimation using the second moment matrix. Harris based detectors tend to extract corner-like structures, while Hessian based detectors tend to find blobs and ridges. DoG [16] is similar to the Hessian detector in the sense that it approximates LoG by the trace of the Hessian matrix. DoG tends to extract points at isotropic blob structures.

Some recent feature detectors also based on scale space are FAST [19, 20] and its variants: AGAST [21], ORB [22], BRISK [23], and FREAK [24]. The FAST method in [19, 20] proposed an efficient approach for corner detection, based on the comparison of pixel values on a ring centered at a feature point. The AGAST [21] detector is an improved version for accelerated performance of FAST. The recent ORB [22] detector is a rotation-invariant extension of FAST. The detector used in BRISK [23] is a multi-scale AGAST. It searches for

maxima in scale-space using the FAST score as a measure of saliency, and estimates the scale of each keypoint in the continuous scale-space. The FREAK [24] method uses the same detector as the one used in BRISK.

The second type of feature detectors is based on tree-based shape space. Although its original definition is quite different, the Maximally Stable Extremal Regions (MSER) [1] are easily understandable using Min-tree and Max-Tree. Indeed, as shown in [25], the MSER algorithm extracts the regions (nodes) that correspond to local minima of a stability function along the path to the root of the tree. The stability function  $\mathcal{A}_q$  of a given node  $\mathcal{N}$  is given by the difference between the area of some (grand-)parent  $\mathcal{N}_\Delta^+$  and some (grand-)child  $\mathcal{N}_\Delta^-$ , divided by the area of the node itself. It is given by:

$$\mathcal{A}_q(\mathcal{N}) = (|\mathcal{N}_\Delta^+| - |\mathcal{N}_\Delta^-|)/|\mathcal{N}|, \quad (1)$$

where  $|\cdot|$  denotes the cardinality,  $\mathcal{N}_\Delta^+$  and  $\mathcal{N}_\Delta^-$  are respectively the lowest ancestor and the highest descendant such that  $|f(\mathcal{N}_\Delta^+) - f(\mathcal{N})| \geq \Delta$  and  $|f(\mathcal{N}) - f(\mathcal{N}_\Delta^-)| \geq \Delta$ , and  $\Delta$  is a stability range parameter that fixes the intensity level difference. It is reported in [6] that MSER achieves state-of-the-art repeatabilities and regions accuracies. It is also very efficient. However, the number of detected features are comparatively small, which limits its ability for some applications like image registration and 3D reconstruction. Perdoch *et al.* [26] propose the Stable Affine Frame (SAF), for which only local stability is required. Many more features are obtained with a comparable repeatability score. However, it is much slower than MSER.

A complete review of invariant feature detectors is out of scope of this paper. The interested reader is referred to Tuytelaars and Mikolajczyk's survey [27]. Although these detectors are widely used in many applications in the computer-vision community, direct application of these detectors to other fields can be difficult. For example, surgical navigation has shown significant difficulties due to the free-form tissue deformation and changes of visual appearance of surgical scenes, and alternate solutions have been proposed [28].

## III. TREE-BASED MORSE REGIONS

In this section, we describe our proposed topology-based local invariant features detector called the *Tree-Based Morse Regions (TBMR)*. The TBMRs are extracted from the shape space built from the image using the Max-tree  $\mathcal{T}_M$  and Min-tree  $\mathcal{T}_m$ . In Section III-A, we briefly review the Morse theory [11, 12] and the original notion of space of shapes. In Section III-B, we show how to extract “interest” regions from these shape spaces. The algorithm of TBMRs extraction based on component trees is described in Section III-C.

### A. Tree-based Morse Theory and Shape Spaces

The aim of Morse theory is to describe the topological changes of the (iso)level sets of a real-valued function in terms of its critical points. Recall that a Morse function is a smooth function  $f$  whose critical points (*i.e.*, points where  $\nabla f = 0$ ) are isolated. Critical points are minima, maxima, and saddle points of  $f$ . The topology of  $f$  is directly linked to the analysis of those critical points.

The use of Morse theory is not new in computer vision: see, *e.g.* the use of the contour tree [29, 30] and the Reeb graph [31, 32] for shape matching. However, the Morse function is not an adequate model for an image, as it prevents the existence of plateaux for example. A consequence is that we want to deal with *regional* extrema and saddle *regions* instead of isolated points. The chapter 4 of [12] describes several notions of critical values and prove that they are equivalent to the critical points of Morse theory. In  $\mathbb{R}^2$ , the Morse structure of an image is extracted from its upper and lower level sets, and is equivalent to the notion of critical value used in Morse theory.

For any  $\lambda \in \mathbb{R}$ , the upper level sets  $\mathcal{X}_\lambda$  and lower level sets  $\mathcal{X}^\lambda$  of an image  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  are defined respectively by  $\mathcal{X}_\lambda(f) = \{p \in \mathbb{R}^2 \mid f(p) \geq \lambda\}$  and  $\mathcal{X}^\lambda(f) = \{p \in \mathbb{R}^2 \mid f(p) \leq \lambda\}$ . Thanks to the inclusion relationship, the connected components of upper (*resp.* lower) level sets can be organized into a tree structure, that is called the Max-tree  $\mathcal{T}_M$  (*resp.* Min-tree  $\mathcal{T}_m$ ).

Let us intuitively describe the building of the Max-tree along with the topological changes in the level sets. Imagine that the surface is completely covered by water, and that the level of water slowly decreases. Islands (regional maxima) appear. These islands form the leafs of the Max-tree. As the level of water decreases further, islands grow, building the branches of the Max-tree. Sometimes, at a given level, several islands merge into one connected piece. Such pieces are the forks of the Max-tree, *i.e.*, the nodes of the tree with several children. We stop when all the water has disappeared. The emerged area forms a unique component: the root of the tree representing the whole image.

Following this intuitive description, the critical regions we consider in this paper are the extrema of  $f$ , corresponding to the leaves of the Max-tree  $\mathcal{T}_M$  and of the Min-tree  $\mathcal{T}_m$  of  $f$ , and the saddle points of  $f$ , corresponding to the nodes of these trees having several children.

According to Morse theory, critical points provide essential geometric information. This would lead us to consider extrema and saddle points. Indeed, these are invariant under any increasing contrast change. Unfortunately, critical points are difficult to use directly: (i) many of them are due to noise, and (ii) there is a crucial absence of information about their scale. Typically, most extrema are single pixels, which do not provide a clue about an adapted neighborhood to use as descriptor. Notice that applying a grain filter to the image does not solve the problem. Most regional extrema will then have the size of the grain filter, which is an artifact of the method. For this reason, and as we detail in the next section III-B, we are led to use the largest regions containing the critical region without a topological change. In other words, we are associating scales to critical regions, corresponding to the components containing them just before their merge. However, in this paper, the exact positions of the critical points are ignored, since at the end the centroids of the components are used as feature points.

An example of a synthetic image, together with its Max-tree and its Min-tree, is shown in Fig. 2. In this figure, critical regions are highlighted with a red circle.

Such types of tree-based image representations feature sev-

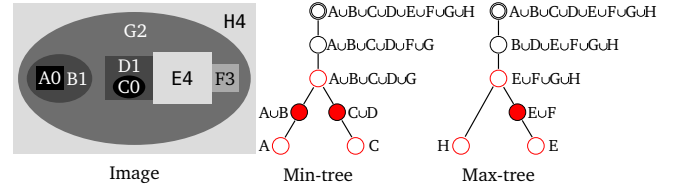


Fig. 2. A synthetic image and the corresponding Min-tree (middle) and Max-tree (right) representation. The critical regions are represented by red circles: (1) nodes having more than one child, and (2) leaf-nodes. The filled regions are the TBMRs.

eral interesting properties. Specifically, they are invariant to affine contrast changes and covariant to continuous (topological) transformations. Furthermore, efficient algorithms with linear or quasi-linear complexity are available [13, 14, 15].

As described in section II, many feature detection methods are based on a scale-space. The *causality principle* is certainly the most fundamental principle of multi-scale analysis [33]. From this principle, for any couple of scales  $\lambda_2 > \lambda_1$ , the “structures” found at scale  $\lambda_2$  should find a “cause” at scale  $\lambda_1$ . Both  $\mathcal{T}_M$  and  $\mathcal{T}_m$  share such a property: indeed, an extremum that appears at a given level  $\lambda_1$  gives rise to a whole branch of the tree, which corresponds to a region from very fine to the whole image. Thus, both  $\mathcal{T}_M$  and  $\mathcal{T}_m$  can be seen as a multi-scale image representation, that we call a shape-space. Such shape-spaces have another interesting property: contrary to scale-space, the contours of a given shape (connected component) correspond to actual contours in the image, without any “blurring” due to convolution with a kernel. These properties of the tree-based image representations make the shape space very appropriate for local invariant feature detection. Indeed, these interesting properties of the shape-space made from the connected components of both upper and lower level sets, coupled with an attractive criterion for regions selection (see Section II), are the reasons of the success of MSER. However, the MSER criterion, based on Eq. (1), is not invariant to contrast changes (contrarily to the TBMR topological criterion detailed in the next section).

### B. Feature extraction based on Morse theory and shape space

The next stage is to associate a scale to each critical region. A critical region corresponds to a change of topology in its tree: either an apparition (leaves corresponding to extrema) or a merge (region resulting from a fork). Thus, on a branch of the tree between two critical regions, there is no topological change in the tree. In other words, a region that is not critical is topologically equivalent to the first critical region we encounter going from the region to the leaves of the tree. Conversely, a critical region  $\mathcal{R}_c$  is topologically equivalent to any region that contains  $\mathcal{R}_c$  but no other disjoint critical region. As we want as much context as reasonable to encode the region, a good “scale” choice for representing a critical region is the largest region to which it is topologically equivalent in its tree. We call such a region a *Tree-based Morse Region* (TBMR).

In practice, we do not consider the TBMRs that are either too small or too big. Discarding small regions is performed

```

TBMR_Extraction( $f$ ,  $\lambda_{\min}$ ,  $\lambda_{\max}$ )
 $\mathcal{T}_M \leftarrow \text{Compute\_MaxTree}(f)$ 
 $\mathcal{S} \leftarrow \text{TBMR\_Tree\_Extraction}(\mathcal{T}_M, \lambda_{\min}, \lambda_{\max})$ 
 $\mathcal{T}_m \leftarrow \text{Compute\_MinTree}(f)$ 
 $\mathcal{S} \leftarrow \mathcal{S} \cup \text{TBMR\_Tree\_Extraction}(\mathcal{T}_m, \lambda_{\min}, \lambda_{\max})$ 
return  $\mathcal{S}$ 

TBMR_Tree_Extraction( $\mathcal{T}$ ,  $\lambda_{\min}$ ,  $\lambda_{\max}$ )
 $\mathcal{S} \leftarrow \emptyset$ 
foreach  $\mathcal{N}$  in  $\mathcal{T}$  do  $\text{numChildren}(\mathcal{N}) \leftarrow 0$ 
foreach  $\mathcal{N}$  in  $\mathcal{T}$  do
  | if  $\text{area}(\mathcal{N}) \geq \lambda_{\min}$  then
  | |  $\text{++numChildren}(\text{parent}(\mathcal{N}))$ 
foreach  $\mathcal{N}$  in  $\mathcal{T}$  do
  | if  $\text{area}(\mathcal{N}) < \lambda_{\max}$  and  $\text{numChildren}(\mathcal{N}) = 1$  and
  | |  $\text{numChildren}(\text{parent}(\mathcal{N})) \geq 2$  then
  | |  $\mathcal{S}.\text{insert}(\mathcal{N})$ 
return  $\mathcal{S}$ 

```

**Algorithm 1:** Extraction of TBMRs. Too small regions (area less than  $\lambda_{\min}$ ) do not contribute to topological changes, and too large regions (area greater than  $\lambda_{\max}$ ) are discarded.

before analyzing the tree structure, which means that they do not contribute to the topological changes of the tree structure. Discarding them eliminates also some noise without modifying other components. In our experiments, we always set this lower bound at 30 pixels. In addition, regions that meet the image border are considered truncated so we ignore them. In Fig. 2, the TBMRs are drawn with a red disk. The evolution of the number of children of a region, starting from a leaf to the root is illustrated in Fig. 3.

As in most shape-space-based methods, we compute the centroids of the selected regions as the final feature points. The ellipse with the same first and second moments as the detected region is then used as the local patch upon which a descriptor is computed.

### C. Algorithm of TBMRs extraction

An algorithm that extracts the TBMRs is provided in Algorithm 1. It relies on the computation of the two component trees Max-tree and Min-Tree [13, 14, 34, 15, 35]. The TBMRs are extracted from each one of those two trees independently. We first count the number of children of each node, ignoring the nodes that are too small. Then, the TBMRs are the nodes with a unique child and at least one sibling, removing the nodes that are too big. The corresponding regions provide a (kind of) scale associated to a given critical point, and are approximated by an ellipse.

As shown in Algorithm 1, the extraction of TBMRs is composed of (1) shape space building by Min/Max-tree computation and (2) a linear regions selection from the shape space. There exist many efficient algorithms to compute the Min/Max-tree [35]. The union-find-based approach [14] would take  $O(n\alpha(n))$  time for low quantized image (typically 12 bits image or less), where  $n$  is the total number of pixels inside an image  $f$  and  $\alpha$  is a very slow-growing “diagonal inverse” of the Ackermann’s function. Linear algorithms to compute

the Min/Max-tree are also available [13, 15]. The selection of TBMRs from the shape would take  $O(N)$  time, where  $N$  is the total number of nodes of the tree which is no greater than  $n$ . In consequence, the complexity of the TBMRs computation algorithm shown in 1 would be  $O(n\alpha(n))$  or  $O(n)$ .

## IV. COMPARISON WITH SOME RELATED WORK

We focus on the comparison with two detector classes, those based on scale-space and those based on tree-based shape space, notably MSER.

### A. TBMR versus scale-space feature detection

In spirit, TBMR is very similar to those kinds of approaches. TBMR detects critical points (*i.e.*, extrema and saddle points), but does not rely on a scale-space. As described in section III, it uses the space of shapes provided by the Min-tree and Max-tree representations. This space has the main property of scale-space, namely the causality principle [33].

### B. TBMR versus MSER

TBMR can be seen as a variant of MSER, since both rely on Min/Max-tree representations. Indeed, TBMRs are the children of merge nodes (bifurcation). However, TBMRs are never MSERs: indeed, at a merge node, the ratio-of-area criterion proposed in [1] is unstable, and thus the MSER algorithms ignore those nodes.

In practice, the most fundamental difference between MSER and TBMR is related to illumination change, a very common effect in natural images that is reported as an unsolved problem in the literature [8]. Indeed, the MSER stability function shown in Eq. (1) depends on a parameter  $\Delta$  that fixes the intensity level difference of the (grand-)parent and of the (grand-)child actually used for the ratio. That prevents a true invariance of MSER to illumination change. By contrast, TBMR, being purely topological, is truly invariant to affine illumination change. A less fundamental difference concerns the number of parameters of MSER. As TBMR, MSER uses two parameters to remove too large and too small regions. But MSER also requires in the stability function, on top of the parameter  $\Delta$  we just described, a threshold to remove unstable regions, and another parameter to group together detected regions that are similar in terms of position and size. The latter is very important, as there are usually too many local minima of the ratio-of-area criterion, and some of them correspond actually to the “same” (very similar) stable regions. Especially in the case of using a very small intensity level difference, which does not have much sense with respect to the notion of stability. Grouping those similar regions makes the extracted regions less accurate (See the registration result presented in Table II).

Such additional parameters are not needed in TBMR. A last minor difference is also related to the definition of MSER at bifurcations: As we mention above, the stability function is not clearly defined in the presence of bifurcations, *i.e.*, when a node has more than one child. That raises a difficulty in trying to reproduce some results: for example, there exist two

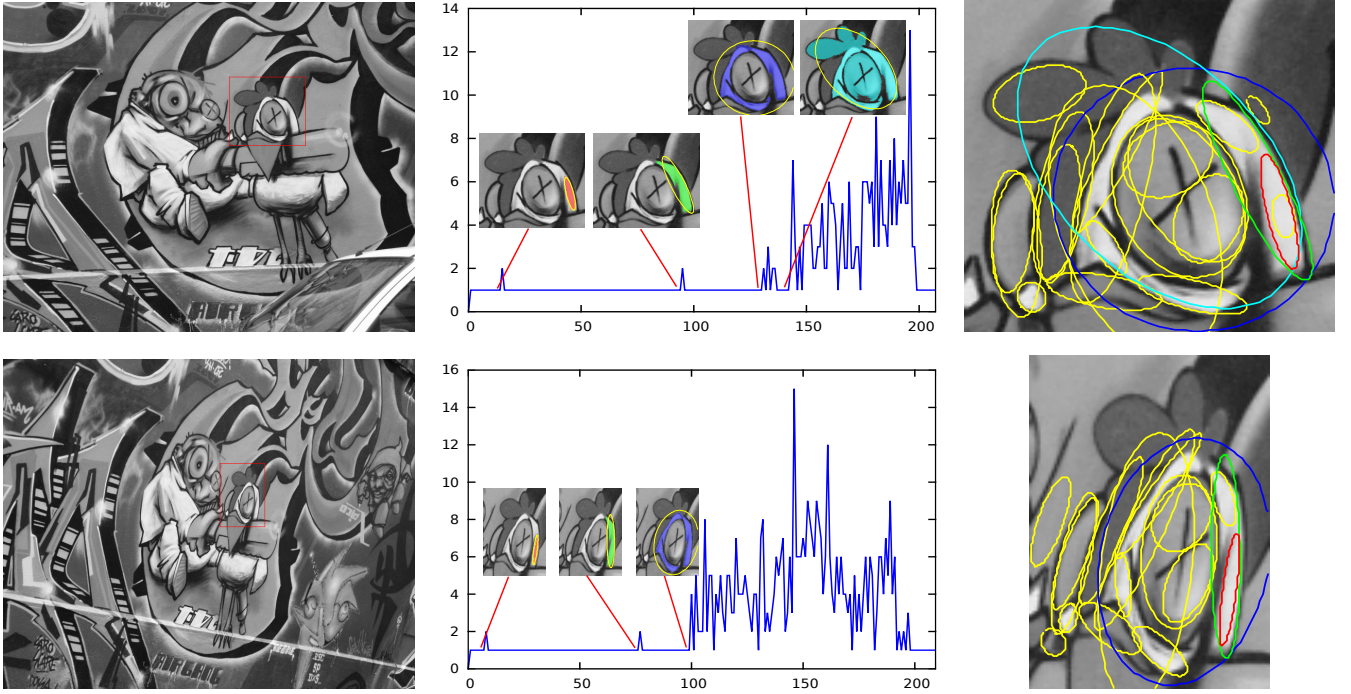


Fig. 3. An example of TBMR extraction. Left: Front view (top) and 30° view (bottom) of “Graffiti” scene [6]. Middle: Evolution of number of children starting from a leaf to the root and some extracted TBMRs (as well as the corresponding ellipses) along the branch; Right: ellipses of the extracted TBMRs inside the corresponding red bounding boxes. Note that only the ellipses of the TBMRs extracted from Max-tree are shown.

public implementations of MSER, one from VLFeat [36], the other from OpenCV, each one using a similar but different stability function. The topological definition of TBMR allows for a perfect reproducibility, whatever the chosen algorithm implementation.

### C. Comparison of time complexities and of running times

We compare the time complexity of TBMR with the similar method MSER, and the widely used DoG detector. Both TBMR and MSER rely on a computation of Min-tree and Max-tree representation whose time complexity is  $O(n\alpha(n))$  or  $O(n)$ . They differ in the extraction of regions from the Min-tree and Max-tree. As described in Section III-C, the TBMR method requires a computation of number of children which has a time complexity of  $O(N)$  and a linear region selection whose time complexity is also  $O(N)$ . The MSER method relies on the stability function given in Eq. (1). For each node, one has to find its (grand-)parent and the (grand-)child with a level difference of  $\Delta$ ; the time complexity of this step is  $O(\Delta N)$ . The next step of the MSER method is to select the local minima of the stability function, this has a time complexity of  $O(N)$ . Consequently, TBMR method and MSER method have a similar time complexity which would be  $O(n\alpha(n))$  or  $O(n)$ . The DoG method is mainly composed of Gaussian blurring, difference of Gaussian, scale-space extrema detection. Since there are just a few number of scales, the time complexity of the naive Gaussian blurring is  $O(\omega^2 n)$ , where  $\omega$  is the window size of the kernel of Gaussian blurring. Linear algorithm for Gaussian blurring is also available thanks to recursive filtering [37, 38]. The difference of Gaussian and scale-space extrema detection both have a  $O(n)$  time

complexity. Finally, the time complexity of DoG is  $O(\omega^2 n)$  or  $O(n)$ .

We have implemented the proposed TBMR method in VLFeat [36]. We compare the running time of TBMR with the public implementation of MSER and DoG already available in VLFeat. The comparison is given in Table I. For each image, each detector is executed 100 times, the average time is shown. All the measurements have been taken on a 3.4GHz/8MB cache Intel core i7-2600, 8GB RAM, Debian 7. The code is compiled with GCC 4.7.2. Note that for the TBMR and MSER methods, the time for region selection is negligible compared to that of the computation of Min-tree and Max-tree. So the two methods have a very similar running time. However, in this public implementation, the accumulated moments for computing the approximated ellipse of each TBMR or MSER region in the Min-tree or Max-tree are computed after the tree construction step. In a more efficient implementation, these moments could have been computed incrementally during the tree construction step. As it is now, the ellipse fitting step takes about 8% of the time of the whole process, depending on the number of moments to compute. TBMR is slower than MSER for the images of *Wall* and *Trees*: indeed, TBMR extracts many more regions than MSER for these two images, hence the ellipse fitting step takes more time. Both TBMR and MSER are faster than DoG. Note that many optimizations are possible for MSER and TBMR. As shown in [39], a more efficient implementation of MSER is  $5.6\times$  faster than DoG. An optimized TBMR should have a similar performance.



Image	Run time [ms.]		
	DoG	MSER	TBMR
Graffiti (800 × 640)	526.6	484.7	471.6
Wall (1000 × 700)	1139.1	1122.9	1124.7
Bark (765 × 512)	579.5	292.1	290.9
Trees (1000 × 700)	1145.4	958.2	1030.4
Leuven (900 × 600)	771.7	524.0	519.6

TABLE I. Averaged computation time for the different detectors applied on several images in the dataset of Mikolajczyk [6]. The corresponding image size is also shown after the name. See text for a discussion.

## V. RESULTS

Qualitative and quantitative comparison of the distribution of TBMR with other popular local feature detectors are illustrated in Section V-A and Section V-B. In section V-C, the repeatability assessment of the TBMR is evaluated using the framework of Mikolajczyk *et al.* [6]. Two applications using local invariant features are presented in sections V-D and V-E to compare TBMR with other widely used detectors. The application to image registration in Section V-D highlights the accuracy and the robustness of TBMR. The experiments are conducted on the Stanford Mobile Visual Search (SMVS) Data Set [40]. In Section V-E, the application to 3D reconstruction using structure from motion is first tested on the dataset of Strecha *et al.* [41], providing the ground truth of the camera positions. The baseline error and angular error measurements reveal the accuracy of TBMR. Then the 3D reconstruction experiments are conducted on sets of images taken in a sunny day around some structure. The structure from motion succeeds in reconstructing a complete 3D model using TBMR, whereas only part of the scenes are reconstructed in 3D model when using other detectors.

In all experiments, the parameters of the corresponding method are set with the recommended values found in the literature. More specifically, for Harris-Affine, Hessian-Affine and MSER, we use the Linux binaries<sup>1</sup> given in [6]. For the Harris-Affine detector and the Hessian-Affine detector, the only parameter is the cornerness threshold which is set to 1000 for Harris-Affine, and to 500 for Hessian-Affine. For the MSER method, the parameter  $\Delta$  in Eq. (1) is set to 10, and the minimum area and maximum area are set to respectively 30 pixels and 1% of the total number of pixels inside the image (we use the same area parameters setting for TBMR). For the DoG method, its public implementation in VLFeat is used. There are five main parameters: the number of octaves is the greatest possible one (*i.e.* roughly  $\log_2(\min(\text{width}, \text{height}))$ ). Each octave is sampled at 3 intermediate scales. The edge threshold is set to 10, which eliminates peaks of the DoG scale space whose curvature is smaller than 10. The peak threshold which filters peaks of the DoG scale space that are too small is set to  $255 \times 0.04/3$ . In the following, except the explicit “DoG octave 0” whose first octave index is set to 0, the starting octave of DoG method is set to  $-1$ . For all the feature matching used in the experiments of image registration in Section V-D and 3D reconstruction in Section V-E, the SIFT descriptors are used. Finally, in all the experiments, the distance ratio of vector angles from the nearest to second

nearest neighbor used in the SIFT descriptors matching is set to 0.6.

### A. Qualitative feature comparison

In this section, we focus on the distribution of the keypoints. This distribution plays an important role for many applications, such as image registration and 3D reconstruction using structure from motion, for which many points covering the object of interest are required.

We compare TBMR with some state-of-the-art local feature detectors (Harris-Affine, Hessian-Affine, DoG octave 0, DoG, and MSER) by visualizing the distribution of the keypoints obtained with each method.

For the shape-space-based MSER and TBMR methods, the centroid of the extracted regions is considered as the detected keypoints. The qualitative comparison is conducted on two images taken against the sunlight in a sunny day, and the keypoints distributions are illustrated respectively in Fig. 4 and 5. We can observe that MSER detects few points, which explains the defects in the 3D reconstruction based on MSER in the example shown in Fig. 1. Harris-Affine, Hessian-Affine, and DoG octave 0 extract a reasonable number of points, but only a few of them are located on the object of interest of the scene; that makes them fail in reconstructing the 3D structure in Fig. 1. Using the default option “octave  $-1$ ” (image of doubled dimensions) for the DoG computation significantly increases the number of detected points. Yet, the additional points are mostly distributed where there were already many points, and not on the object of interest. By contrast, TBMR has a reasonable number of keypoints and they are distributed more uniformly over the whole image. That contributes to its success in obtaining a correct 3D reconstruction; see Fig. 1.

In Fig. 6, we show the qualitative comparison of TBMR with MSER and DoG in the case of a significant change of contrast. MSER and DoG detects very few points in the area with low contrast. By increasing the contrast, MSER and DoG detects some points. On the contrary, TBMRs are perfectly insensitive to contrast change, up to quantization effects. This better performance of TBMR with respect to change of contrast (together with TBMR robustness to viewpoint change, see section V-D below) is one of the main reason for its success in 3D reconstruction; see Fig. 1, Fig. 17, and Fig. 18.

### B. Image coverage evaluation

In order to assess the uniformity of those keypoint distributions obtained with different methods, we first measure the distribution of keypoints along the two image dimensions, as well as the number of extracted points for a set of images taken around some scene objects. In Fig. 7, we show the distribution of keypoints position along the horizontal dimension for the images taken around the objects of scene presented respectively in Fig. 5, Fig. 17, and Fig. 18. These images are taken to make the object of interest presented in the middle (horizontally) of the scene. The distributions shown in Fig. 7 are smoothed by taking the average inside a horizontal window (size is set to 21). The total number of keypoints for each coordinate axes can be obtained by multiplying the distribution

<sup>1</sup>Available on <http://www.robots.ox.ac.uk/~vgg/research/affine>

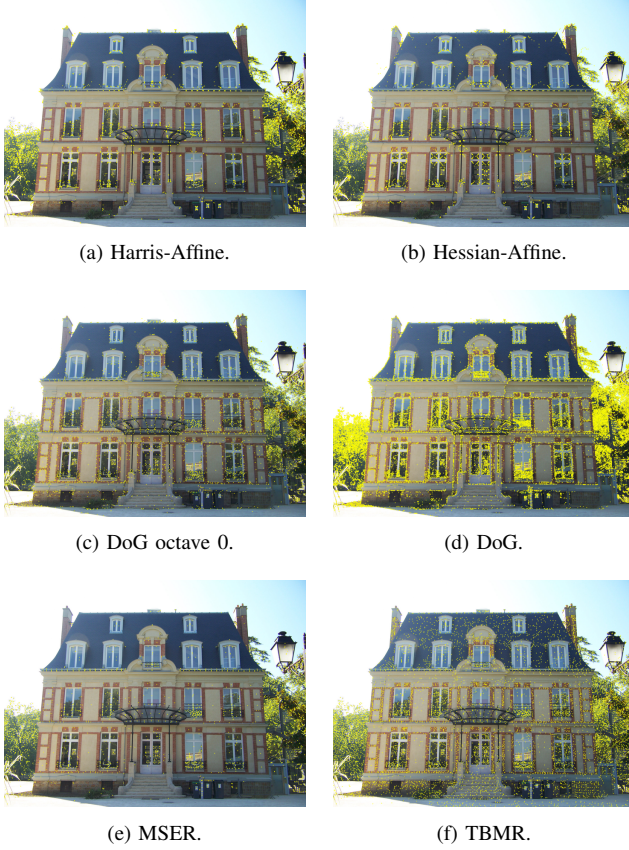


Fig. 4. Qualitative comparison of TBMR with other widely used local feature detectors applied on an image taken against the sunlight (used in Figure 1). Yellow points in the image are the detected keypoints.

score with the average number of keypoints shown in Fig. 7. In order to better visualize the comparison between different methods, only DoG, MSER, and TBMR are shown in this figure. MSER extracts few points; TBMR extracts more points than MSER; DoG has many more points. However, TBMR has the largest part of keypoints that cover the objects of interest in those scenes, which is also one of the reason for its success in Fig. 17 and 18.

We also evaluate how well the keypoints cover the image. First, we propose to dilate the extracted keypoints by a 2D window centered at each point with a certain size (e.g., 31). Then we compute the rate of area covered by the dilated region. Note that for two close keypoints, their dilated regions may have a large part in common, but the common regions count only once. As shown in Fig. 8, MSER covers a small part of the image because of a few extracted points. TBMR covers the image better than the others having a comparable or much larger number of detected points, which confirms the qualitative observation in Section V-A.

### C. Repeatability evaluation

To assess the performance of TBMR, we compare it with some other affine detectors: Harris-Affine and Hessian-Affine, defined on a scale-space, and MSER, defined on the shape space.

We perform the same tests as in [6], and evaluate the



Fig. 5. Qualitative comparison of TBMR with other widely used local feature detectors applied on an image taken against the sunlight. Yellow points in the image are the detected keypoints.

repeatability score based on the overlap error  $\epsilon$ :

$$\epsilon(R_{E_1}, R_{E_2}) = 1 - \frac{R_{E_1} \cap R_{H_{21}^T E_2 H_{21}}}{R_{E_1} \cup R_{H_{21}^T E_2 H_{21}}}, \quad (2)$$

where  $R_E$  represents the elliptic region (i.e., local patch of each extracted feature) defined by  $x^T E x \leq 1$ , and  $H_{21}$  is the ground truth homography between the test and reference image. The repeatability score for a pair of images is then defined as the ratio between the number of region-to-region correspondences established under a certain overlap error (e.g., 40% is used in this paper) and the smaller number of regions in the compared images. Another evaluated measurement is the absolute number of correspondences. A high repeatability score and a large number of correspondences are normally desired.

Some results applied to the sequences *Wall* (viewpoint change), *Bark* (scale change), *Trees* (blur), *Leuven* (light change) [6] are illustrated in Fig. 9. Compared to the scale-space-based approaches (i.e., Harris-Affine and Hessian-Affine), the TBMR achieves a competitive repeatability score and a significantly higher number of correspondences, except for the blur sequence *Trees*. The explanation is that the topology of the image is damaged by the blur. For the same reason, the performance is poor on the UBC sequence (not shown here) dedicated for testing robustness to strong JPEG compression artifacts. Such defects (blur, JPEG artifacts) are better handled by the scale-space-based methods. Compared

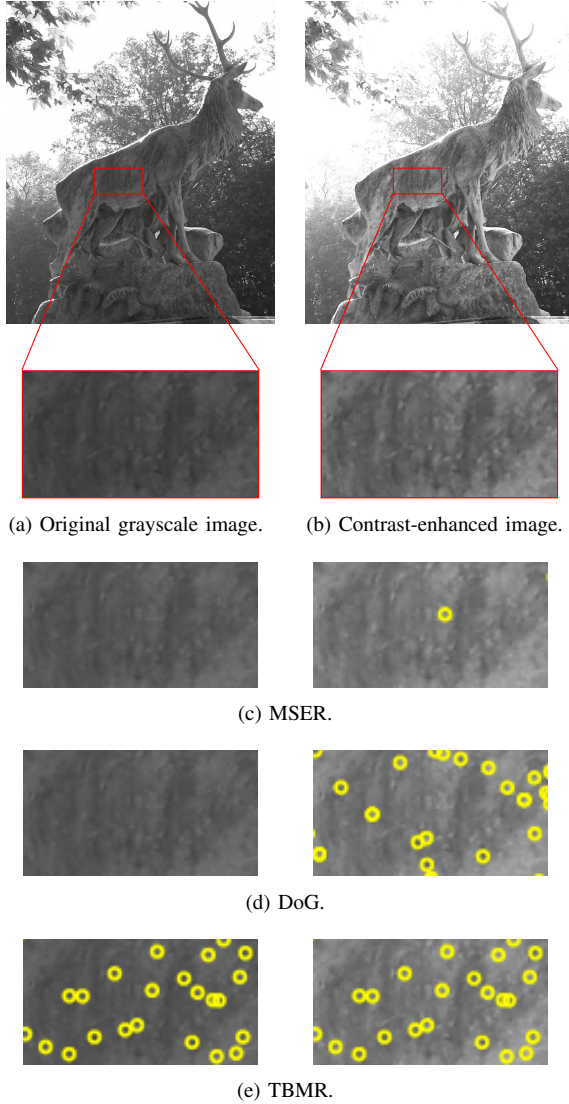


Fig. 6. Qualitative comparison of TBMR with MSER and DoG applied on images with increasing change of contrast. Yellow points in the image are the detected keypoints. Left: Keypoints extracted from the zoomed part of the original image (a); Right: Keypoints extracted from the zoomed part of the contrast-enhanced image (b).

with MSER, which is also based on the shape space, TBMR has a comparable repeatability score, but a significantly higher number of correspondences thanks to the contrast independent property of TBMR. More extensive tests can be found in the supplementary material accompanying this article.

Some extra experiments on other datasets (such as the dataset of DTU [8]) that contains more images, will be considered as a future work.

#### D. Image registration

Image registration methods use the local features to establish a correspondence between a number of interest points (e.g., the centroids of the detected elliptical regions) in images. These one-to-one correspondences are then used to estimate the transformation, thereby establishing point-to-point correspondence between the reference image and the target image.

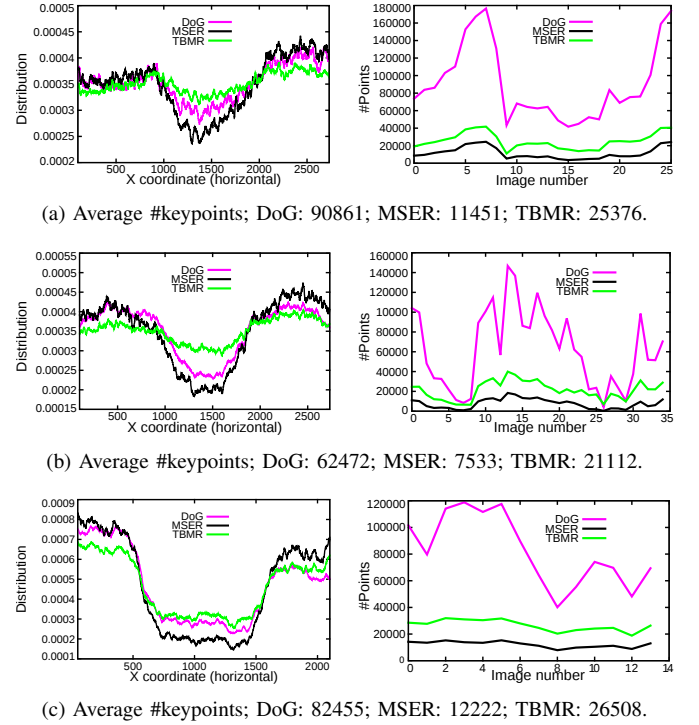


Fig. 7. Horizontal distribution of the keypoints (left) and number of extracted keypoints (right), for the multi-view images taken around the objects of scene presented in respectively Fig. 5, Fig. 17, and Fig. 18 (top to bottom).

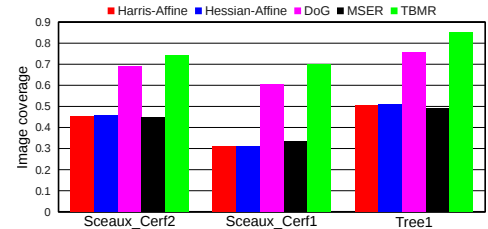


Fig. 8. Image coverage evaluation by dilating a keypoint with a 2D window having size  $31 \times 31$ ; Vertical axis represents the average rate of area covered by the dilated region for multi-view images presented in Fig. 5, Fig. 17, and Fig. 18.

Hence, the accuracy and robustness of the local features are crucial to the quality of image registration results. There should also be enough pairs of matched points so that the estimation of the parameters of the transformation between images is possible.

For these experiments, we use the work of Moisan *et al.* [42]. It is based on the Optimized Random Sampling Algorithm (ORSA) proposed by Moisan and Stival [43], a variant of RANSAC algorithm introducing an *a contrario* criterion [44] to avoid to set thresholds for inlier/outlier discrimination. ORSA is used to estimate the homography registering both images. When the homography is assessed, a panorama is built by stitching the images in the coordinate frame of the second image.

1) *Quantitative benchmark on dataset of Mikolajczyk [6]:* We first benchmark TBMR with MSER, Harris-Affine, Hessian-Affine, and DoG on the public dataset of Mikolajczyk *et al.*, where the ground truth of homography  $H$  is available.



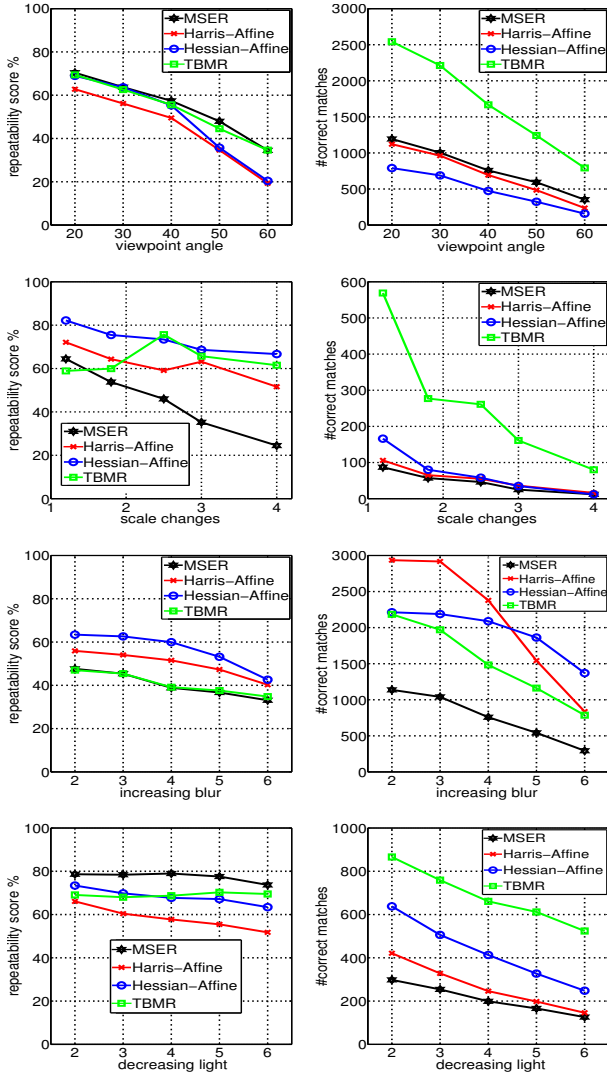


Fig. 9. Repeatability score (left) and number of correspondences (right) for the sequences *Wall*, *Bark*, *Trees*, *Leuven* (top to down).

For each sequence of images in this dataset, the homography between the first image and the others are provided. We have conducted the image registration experiments on each sequence of images in this dataset. The first image is used as reference image, and the other images with different levels of transformation (*e.g.* viewpoint change, scale change, etc.) are considered as target images.

For each point  $p$  in the reference image (the first image in each sequence) which belongs to the inlier matched points used to estimate the homography  $H'$ , we measure the euclidean distance between the points  $H \times p$  and  $H' \times p$  in the registered image. This distance measures the registration errors: the smaller the distance, the more precise the registration.

Distributions of distances in image registration by homography on “Graffiti” (viewpoint changes) image pairs are shown in Fig. 10 (a - e). The average of distances for each image pair is shown in Fig. 10 (f). In general, TBMR performs better than the others.

Method	Canon	Palm
MSER	65	48
MSER2	93	81
Harris-Affine	92	91
Hessian-Affine	97	88
DoG	98	97
TBMR	97	91

TABLE II. Benchmark of image registration results using different local feature detectors on images taken with “Canon” and “Palm” cameras in SMVS dataset; MSER2 represents the MSER method with the margin value set to 2.

2) *Qualitative results*: The experiments are conducted on the CD-covers of the Stanford Mobile Visual Search (SMVS) Data Set [40]. As the CD-covers are planar scenes, they are amenable to homography registration. We experimented Harris-Affine, Hessian-Affine, MSER, DoG, and the proposed TBMR to compute point correspondences between images using the SIFT descriptors of Lowe [4]. The point correspondences are the input of ORSA.

For the images of Fig. 11 and Fig. 13, Harris-Affine, Hessian-Affine and MSER all fail in estimating the homography due to insufficient number of correspondences. DoG also fails for images in Fig. 13. In Fig. 11, although DoG achieves a homography, the registration is inaccurate at the top left corner, whereas the TBMR results in a meaningful homography in all cases. A chessboard mix of the two registered images for Fig. 11 and Fig. 13 are given respectively in Fig. 12 and in Fig. 14, from which the qualitative results can be better visualized. Note that for these image registration examples, there is no ground truth for these images and finding a relevant metric for such poor quality images is a challenge in itself, so only visual inspection is left to the reader’s appreciation.

There are 100 different CD-covers in the tested dataset. We have conducted the image registration experiments on the images taken respectively with “Canon” and “Palm” cameras, and the images of reference. The images taken with “Canon” are less blurred than the ones with “Palm”. In Table II, we show the number of images having a registration result, obtained using different local feature detectors. MSER with standard parameters (margin value  $\Delta$  is set to 10) fails for many images; By lowering considerably the margin value ( $\Delta = 2$ ), the performance gets better; In general, TBMR performs better than MSER (with the standard parameters and a very small margin value), Harris-Affine, and Hessian-Affine. TBMR is on par with DoG when the blur is not very important. In addition, among the tested images, we observe that for those images having registration results, the ones using TBMR are more precise in location than those compared methods. But when the blur is important, DoG is better, because topology of the image is damaged by the blur. However, as shown in [45], the multi-resolution detection improves the performance of MSER under blur. We would expect the same improvements by applying a multi-resolution analysis.

### E. 3D reconstruction

Structure from Motion (SfM) is a popular process of estimating a three-dimensional structure from a sequence of two-dimensional images. The SfM algorithms take multi-



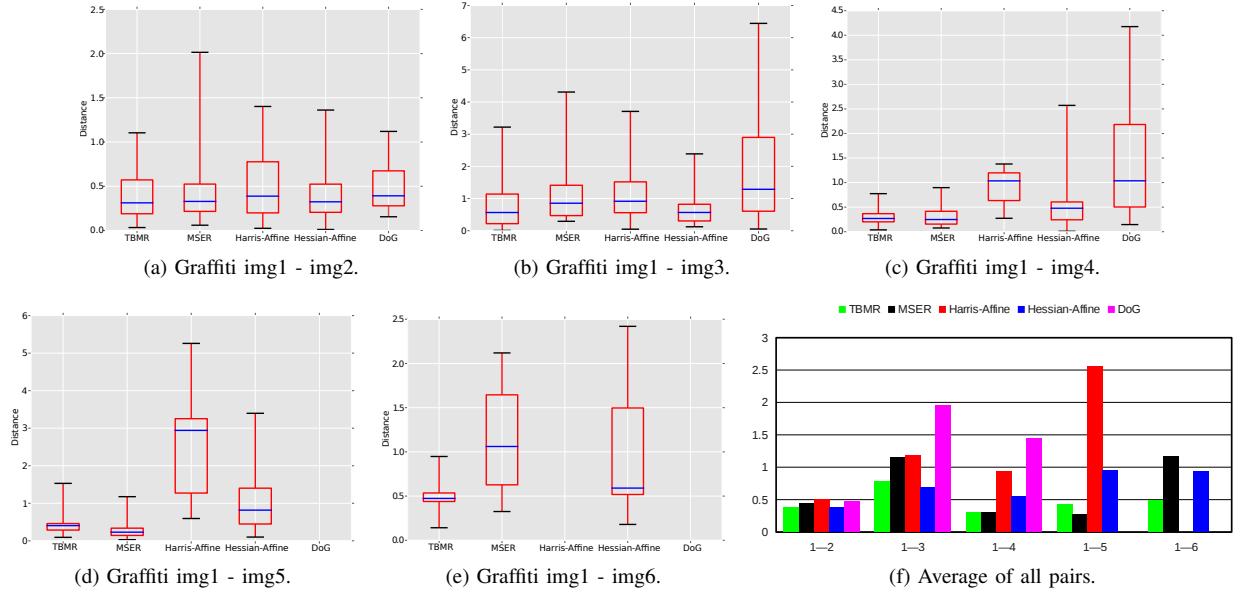


Fig. 10. Distributions of errors in image registration by homography on “Graffiti” image pairs, which comes with a ground truth. DoG fails registering pairs (1,5) and (1,6); Harris-Affine fails on (1,6). There is a slight advantage for MSER on pair (1,5), while on all other pairs TBMR performs better.



Fig. 11. Homographic registration of a pair of images. The result obtained using DoG (b) is not as accurate as the one based on TBMR (c); see the zoomed top left corner. Harris-Affine, Hessian-Affine, and MSER all fail in registering this pair of images.



Fig. 12. Chessboard mix of the two registered images using respectively DoG for (a) and TBMR for (b).

view stereo images along with the internal camera calibration information as input, and yield a sparse 3D point cloud, camera orientations and poses in a common 3D coordinate system. The feature points are used in the phase of model estimations, including homography, fundamental and essential matrices, and camera poses. These estimations are crucial to the quality of 3D reconstruction. Therefore, the accuracy and robustness of the local invariant feature detectors is a defining aspect of



Fig. 13. Homographic registration of a pair of images. Left: reference image. Middle: target image. Right: registration result using TBMR. None other tested detector yields a correct registration.

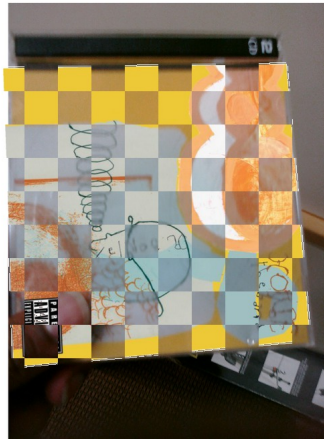


Fig. 14. Chessboard mix of the two registered images using TBMR.

the 3D reconstruction result.

To assess the performance of the feature detectors in this application, we use the software of Moulon *et al.* [10] which relies on an *a contrario* criterion [44] instead of setting thresholds for model estimation in SfM pipelines. In other respects, their pipeline is similar to the one of the popular Bundler software [46]: some initial pair of images are selected and a two-view 3D reconstruction is performed; the other images are then sequentially added, each time refining the 3D scene with bundle adjustment, an iterative optimization method based on Levenberg-Marquardt algorithm. It is shown in their work that the adaptive SfM outperforms the state-of-the-art methods with significant precision improvements. The pipeline may stop prematurely if not enough point correspondences are found when adding extra images. The 3D reconstruction results in the following are obtained thanks to their software OpenMVG [47], an open source software package for SfM.

1) *Quantitative benchmark*: We first benchmark TBMR with Harris-Affine, Hessian-Affine, MSER, and DoG on the public dataset of Strecha *et al.* [41], where the ground truth of the camera orientations and poses are available. The SIFT descriptors of Lowe [4] are again used to establish for each detector the one-to-one correspondences between the feature points. The quality of the 3D reconstruction is tested in terms of the precision of estimated camera orientations and poses. The baseline errors and angular errors compared to the ground truth are illustrated in Fig. 15. For each sequence, the absence of the curve corresponding to some detectors means that it fails to calibrate all the cameras, or that the baseline and angular

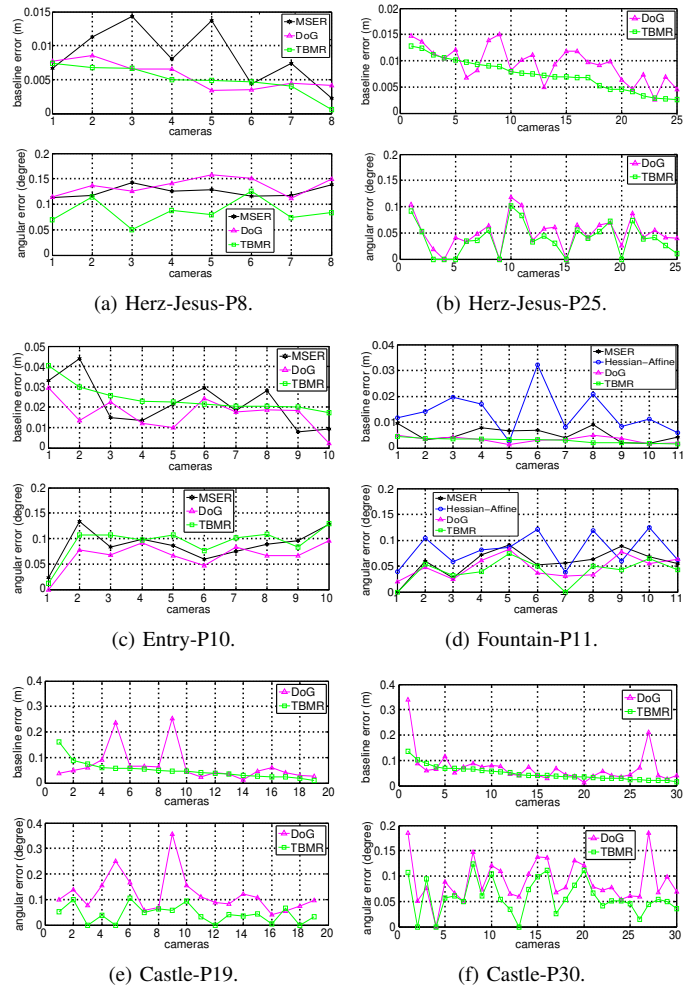


Fig. 15. Evaluation of camera calibration based on baseline error and angular error applied on the dataset in [41].

errors are too high compared with others. Harris-Affine fails or the measurements are too high for all the sequences. Hessian-Affine fails or the measurements are too high for five sequences, while MSER works for three sequences. DoG and TBMR succeed for all the sequences. Compared to Harris-affine, Hessian-affine, and MSER, TBMR is more robust and behaves better in terms of baseline errors and angular errors. Compared to DoG, in most cases (especially for the sequences of *Herz-Jesus-P8* and *Castle-P19*), TBMR performs better based on the baseline and angular errors. In Fig. 16, we show also the recall rate of good tracks, the ratio between the amount of final maintained tracks, and the number of input tracks found between the feature points. The absolute number of final maintained tracks is presented as well. The highest recall rate and the largest number of 3D points obtained with TBMR also reveal its robustness.

2) *Qualitative results*: We have also tested the SfM method with different detectors on some other set of images taken in a sunny day. Since the SfM produces a sparse 3D point cloud, not a dense 3D reconstruction, the PMVS software [9] is used to densify the 3D points. Note that PMVS relies on the interest points of DoG and on the Harris corners to reconstruct the 3D structures based on the estimated model. Consequently,

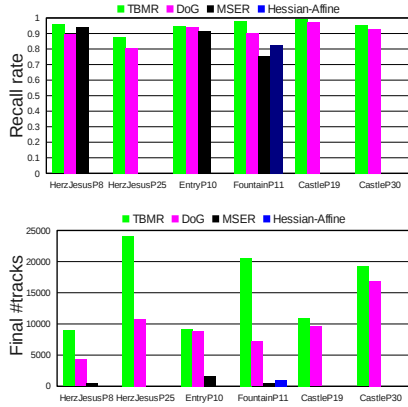


Fig. 16. Recall rate (top) of the tracks and the absolute number (bottom) of final maintained tracks used to yield the sparse 3D points.

this may result in a lack of 3D points in a region (*e.g.* the surface near the bottom of the 3D structures on the right side of Fig. 18 using the TBMR). By integrating the TBMR in PMVS, we would expect the reconstructed 3D structures to be denser thanks to its contrast independent property. In Fig. 17, the DoG detector fails to reconstruct the half of the scene where shadows are present. In Fig. 18, the DoG detector behaves similarly, whereas TBMR works well in all cases: the complete 3D structures are reconstructed. Note that Harris-Affine, Hessian-Affine, and MSER perform even worse than DoG, so the corresponding results are not presented in Fig. 17 and 18. Those examples also confirm the importance of the invariance to illumination changes, as pointed out by Aanæs *et al.* [8]. The resulting 3D reconstruction can be better appreciated in the video in the supplementary material accompanying this article.

## VI. CONCLUSION

We have introduced a topological approach to local feature detection motivated by Morse theory. It relies on the critical points (*i.e.*, minima, maxima, saddle points) of the image and on the shape space given by the Max-tree and the Min-tree built from the image. More precisely, we use the critical regions that are the leaves and nodes with bifurcation in the Max- and Min-trees. To each critical region, we propose to associate the largest region from the shape space that contains it but does not contain any disjoint critical region. We have shown that the proposed method, called TBMR, is truly contrast independent and almost parameter-free. Besides, TBMR is fast to compute with a linear or quasi-linear complexity.

Experimentally, we showed on standard data [6] that the developed TBMR achieves a reasonable repeatability and a significantly higher number of features extracted in images. We have also conducted some experiments on two applications relying on local features using public data sets. The homographic registration results and 3D reconstruction results demonstrated the accuracy and robustness of TBMR compared to other state-of-the-art detectors. Its contrast-invariance property and its robustness to viewpoint change make TBMR a method of choice in numerous practical situations.

In the future, we plan to explore the choice of the associated region from the shape space for each critical point. In the definition of TBMR, instead of extracting the largest region, we can imagine selecting the most meaningful region based on some significance measure (*e.g.* the average of gradient's magnitude). Yet such a measure should be designed to be invariant to illumination changes and affine transformations. We also would like to introduce an option to control the number of extracted TBMRs through the topological persistence [48]. Another future work is to assess the usefulness of TBMR for applications such as object recognition and tracking.

Supplementary material is available on the IEEE website, on <http://laurentnajman.org/index.php?page=tbmr> and on <http://publis.lrde.epita.fr/xu.14.itip>. It contains some additional experimental results, a video, source code and Linux executables. We will soon make the TBMR available in VLFeat [36], OpenCV and OpenMVG (<http://imagine.enpc.fr/~moulonp/openMVG/>).

## ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their helpful comments that greatly contributed to improve this paper. The authors are also grateful to Pierre Moulon for providing the data and the code for the 3D reconstruction, and to Edwin Carlinet for helping to release the code of TBMR. Part of this work was funded by the Agence Nationale de la Recherche (project STEREO - ANR-12-ASTR-0035 and project KIDICO - ANR-2010-BLAN-0205-03.).

## REFERENCES

- [1] J. Matas, O. Chum, M. Urban, and T. Pajdla, "Robust wide baseline stereo from maximally stable extremal regions," in *Proc. of British Machine Vision Conference*, 2002, pp. 384–393.
- [2] T. Tuytelaars and L. V. Gool, "Matching widely separated views based on affine invariant regions," *International Journal of Computer Vision*, vol. 59, no. 1, pp. 61–85, 2004.
- [3] T. Lindeberg, "Feature detection with automatic scale selection," *International Journal of Computer Vision*, vol. 30, pp. 79–116, 1998.
- [4] D. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, pp. 91–110, 2004.
- [5] C. Schmid, R. Mohr, and C. Bauckhage, "Evaluation of interest point detectors," *International Journal of Computer Vision*, vol. 37, no. 2, pp. 151–172, 2000.
- [6] K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. V. Gool, "A comparison of affine region detectors," *International Journal of Computer Vision*, vol. 65, no. 1-2, pp. 43–72, Nov. 2005.
- [7] P. Moreels and P. Perona, "Evaluation of features detectors and descriptors based on 3d objects," *International Journal of Computer Vision*, vol. 73, no. 3, pp. 263–284, 2007.





(a) 4 among 35 used multi-view images.



(b) 3D result using DoG, 17 images are calibrated.



(c) 3D result using TBMR, 27 images are calibrated.

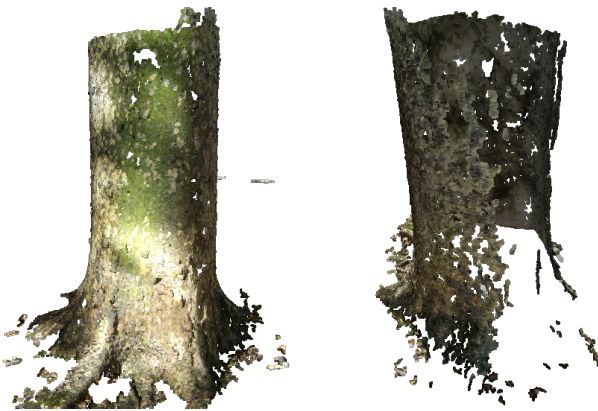
Fig. 17. Densified 3D reconstruction. (a) Some input images. (b) and (c) show the left and right side of the 3D structures reconstructed using respectively DoG and TBMR. The use of DoG misses the right side of the scene; the use of TBMR results in a correct 3D reconstruction; other detectors behave worse than DoG.

- [8] H. Aanæs, A. L. Dahl, and K. Steenstrup Pedersen, "Interesting interest points," *International Journal of Computer Vision*, vol. 97, no. 1, pp. 18–35, 2012.
- [9] Y. Furukawa and J. Ponce, "Accurate, dense, and robust multiview stereopsis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 8, pp. 1362–1376, 2010.
- [10] P. Moulon, P. Monasse, and R. Marlet, "Adaptive structure from motion with a contrario model estimations," in *Proc. of Asian Conference on Computer Vision*, 2012.
- [11] J. Milnor, *Morse Theory*, ser. Annals of Mathematics Studies. Princeton University Press, 1963, vol. 51.
- [12] V. Caselles and P. Monasse, *Geometric Description of Images as Topographic Maps*, 1st ed. Springer Publishing Company, Incorporated, 2009.
- [13] P. Salembier, A. Oliveras, and L. Garrido, "Antiextensive connected operators for image and sequence processing," *IEEE Transactions on Image Processing*, vol. 7, no. 4, pp. 555–570, 1998.
- [14] L. Najman and M. Couprie, "Building the component tree in quasi-linear time," *IEEE Transactions on Image Processing*, vol. 15, no. 11, pp. 3531–3539, 2006.
- [15] D. Nistér and H. Stewénus, "Linear time maximally stable extremal regions," in *Proc. of European Conference on Computer Vision*, 2008, pp. 183–196.
- [16] D. Lowe, "Object recognition from local scale-invariant features," in *Proc. of IEEE International Conference on Computer Vision*, 1999, pp. 1150–1157.
- [17] C. Harris and M. Stephens, "A combined corner and edge detector," in *Proc. of Fourth Alvey Vision Conference*, 1988, pp. 147–151.
- [18] K. Mikolajczyk and C. Schmid, "Scale & affine invariant interest point detectors," *International Journal of Computer Vision*, vol. 60, no. 1, pp. 63–86, 2004.
- [19] E. Rosten and T. Drummond, "Machine learning for high-speed corner detection," in *Proc. of European Conference on Computer Vision*, 2006, pp. 430–443.
- [20] E. Rosten, R. Porter, and T. Drummond, "Faster and better: A machine learning approach to corner detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 1, pp. 105–119, 2010.
- [21] E. Mair, G. D. Hager, D. Burschka, M. Suppa, and G. Hirzinger, "Adaptive and generic corner detection based on the accelerated segment test," in *Proc. of European Conference on Computer Vision*, 2010, pp. 183–196.
- [22] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *Proc. of IEEE International Conference on Computer Vision*, 2011, pp. 2564–2571.
- [23] S. Leutenegger, M. Chli, and R. Y. Siegwart, "Brisk: Binary robust invariant scalable keypoints," in *Proc. of IEEE International Conference on Computer Vision*, 2011, pp. 2548–2555.
- [24] A. Alahi, R. Ortiz, and P. Vandergheynst, "Freak: Fast retina keypoint," in *Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2012, pp. 510–517.

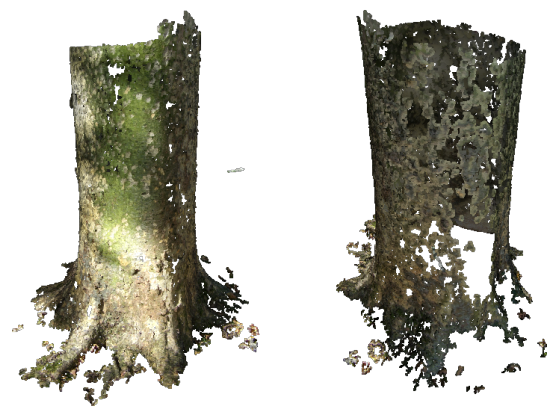




(a) 4 among 14 used multi-view images.



(b) 3D result using DoG, 13 images are calibrated.



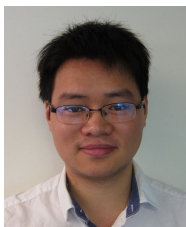
(c) 3D result using TBMR, 14 images are calibrated.

Fig. 18. Densified 3D reconstruction. (a) Some input images. (b) and (c) show the left and right side of the 3D structures reconstructed using respectively DoG and TBMR. The use of DoG misses the right side of the scene; the use of TBMR results in a correct 3D reconstruction; other detectors behave worse than DoG.

- [25] M. Donoser and H. Bischof, “Efficient maximally stable extremal region (MSER) tracking,” in *Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2006, pp. 553–560.
- [26] M. Perdoch, J. Matas, and S. Obdrzlek, “Stable affine frames on isophotes,” in *Proc. of IEEE International Conference on Computer Vision*, 2007, pp. 1–8.
- [27] T. Tuytelaars and K. Mikolajczyk, *Local Invariant Feature Detectors: A Survey*. Hanover, MA, USA: Now Publishers Inc., 2008.
- [28] S. Giannarou, M. Visentini-Scarzanella, and G.-Z. Yang, “Probabilistic tracking of affine-invariant anisotropic regions,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 130–143, 2013.
- [29] I. S. Kweon and T. Kanade, “Extracting topographic terrain features from elevation maps,” *Journal of Computer Vision, Graphics and Image Processing: Image Understanding*, vol. 59, no. 2, pp. 171–182, 1994.
- [30] M. Van Kreveld, R. van Oostrum, C. Bajaj, V. Pascucci, and D. Schikore, “Contour trees and small seed sets for isosurface traversal,” in *Proceedings of the Thirteenth Annual Symposium on Computational Geometry*, 1997, pp. 212–220.
- [31] G. Reeb, “Sur les Points Singuliers d’une Forme de Pfaff Complètement Intégrable ou d’une Fonction Numérique,” *Comptes Rendus Acad. Sciences*, vol. 222, pp. 847–849, 1946.
- [32] S. Takahashi, T. Ikeda, Y. Shinagawa, T. L. Kunii, and M. Ueda, “Algorithms for extracting correct critical points and constructing topological graphs from discrete geographical elevation data,” *Comput. Graph. Forum*, vol. 14, no. 3, pp. 181–192, 1995.
- [33] J. J. Koenderink, “The structure of images,” *Biological Cybernetics*, vol. 50, no. 5, pp. 363–370, 1984.
- [34] C. Berger, T. Géraud, R. Levillain, N. Widynski, A. Bailard, and E. Bertin, “Effective component tree computation with application to pattern recognition in astronomical imaging,” in *Proc. of IEEE International Conference on Image Processing*, vol. 4, 2007, pp. 41–44.
- [35] E. Carlinet and T. Géraud, “A comparison of many max-tree computation algorithms,” in *Mathematical Morphology and Its Applications to Signal and Image Processing—Proc. of the International Symposium on Mathematical Morphology (ISMM)*, ser. Lecture Notes on Computer Sciences, vol. 7883. Springer, 2013, pp. 73–85.
- [36] A. Vedaldi and B. Fulkerson, “VLFeat: An open and portable library of computer vision algorithms,” 2008,

<http://www.vlfeat.org/>.

- [37] R. Deriche, “Recursively implementating the gaussian and its derivatives,” INRIA, Tech. Rep., April 1993.
- [38] B. Appleton and H. Talbot, “Recursive filtering of images with symmetric extension,” *Signal Processing*, vol. 85, no. 8, pp. 1546–1556, 2005.
- [39] O. Miksik and K. Mikolajczyk, “Evaluation of local detectors and descriptors for fast feature matching,” in *Proc. of International Conference on Pattern Recognition*, Nov 2012, pp. 2681–2684.
- [40] V. R. Chandrasekhar, D. M. Chen, S. S. Tsai, N.-M. Cheung, H. Chen, G. Takacs, Y. Reznik, R. Vedantham, R. Grzeszczuk, J. Bach, and B. Girod, “The Stanford mobile visual search data set,” in *MMSYS*, 2011, pp. 117–122.
- [41] C. Strecha, W. V. Hansen, L. V. Gool, P. Fua, and U. Thoennessen, “On benchmarking camera calibration and multi-view stereo for high resolution imagery,” in *Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- [42] L. Moisan, P. Moulon, and P. Monasse, “Automatic Homographic Registration of a Pair of Images, with A Contrario Elimination of Outliers,” *IPOL*, 2012, <http://dx.doi.org/10.5201/ipol.2012.mmm-oh>.
- [43] L. Moisan and B. Stival, “A probabilistic criterion to detect rigid point matches between two images and estimate the fundamental matrix,” *International Journal of Computer Vision*, vol. 57, no. 3, pp. 201–218, 2004.
- [44] A. Desolneux, L. Moisan, and J.-M. Morel, “Meaningful alignments,” *International Journal of Computer Vision*, vol. 40, no. 1, pp. 7–23, 2000.
- [45] P.-E. Forssen and D. G. Lowe, “Shape descriptors for maximally stable extremal regions,” in *Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.
- [46] N. Snavely, S. M. Seitz, and R. Szeliski, “Photo tourism: Exploring photo collections in 3D,” *ACM Transactions on Graphics (TOG)*, vol. 25, no. 3, 2006.
- [47] P. Moulon, P. Monasse, and R. Marlet, “openMVG: open multiple view geometry,” 2012, <http://imagine.enpc.fr/~moulonp/openMVG/>.
- [48] H. Edelsbrunner, D. Letscher, and A. Zomorodian, “Topological persistence and simplification,” in *Proc. of IEEE Symposium on Foundations of Computer Science*, 2000, pp. 454–463.



**Yongchao Xu** received in 2010 both the engineer degree in electronics & embedded systems at Polytech Paris Sud and the master degree in signal processing & image processing at Université Paris Sud, and the Ph.D. degree in image processing and mathematical morphology at Université Paris Est in 2013. He is currently working at EPITA Research and Development Laboratory (LRDE) and at the Signal and Image Processing Department, Telecom ParisTech as a postdoctoral fellow. His research interests include mathematical morphology, image segmentation, medical image analysis, and local feature detection.



**Pascal Monasse** received a Ph.D. degree in applied mathematics from the University Paris-Dauphine in 2000 and the Habilitation à Diriger les Recherches from École Normale Supérieure de Cachan in 2013. He previously got an engineering degree from École Nationale des Ponts et Chaussées in 1995. He worked as a junior then senior researcher at Cognitech, Inc. in Pasadena, California, on image and video processing from 2001 to 2007. The software product PixL2GPS, of which he was the main developer and project leader while working at Cognitech, received in 2010 the American Technology Awards “The Termans” from the TechAmerica Foundation. He joined the Imagine group at École des Ponts ParisTech as a research scientist in 2008 and is a member of Laboratoire d’Informatique Gaspard Monge (LIGM), Université Paris-Est Marne-la-Vallée. His current research interests include computer vision, stereo reconstruction and mathematical morphology.



**Thierry Géraud** received a Ph.D. degree in signal and image processing from Télécom ParisTech in 1997, and the Habilitation à Diriger les Recherches from Université Paris-Est in 2012. He is one of the main authors of the Olena platform, dedicated to image processing and available as free software under the GPL licence. His research interests include image processing, pattern recognition, software engineering, and object-oriented scientific computing. He is currently working at EPITA Research and Development Laboratory (LRDE), Paris, France.



**Laurent Najman** Laurent Najman received the Habilitation à Diriger les Recherches in 2006 from the University of Marne-la-Valle, a Ph.D. of applied mathematics from Paris-Dauphine University in 1994 with the highest honour (Félicitations du Jury) and an Ingénieur degree from the Ecole des Mines de Paris in 1991. After earning his engineering degree, he worked in the Central Research Laboratories of Thomson-CSF for three years, working on some problems of infrared image segmentation using mathematical morphology. He then joined a start-up company named Animation Science in 1995, as director of research and development. The technology of particle systems for computer graphics and scientific visualisation, developed by the company under his technical leadership received several awards, including the European Information Technology Prize 1997 awarded by the European Commission (Esprit programme) and by the European Council for Applied Science and Engineering and the Hottest Products of the Year 1996 awarded by the Computer Graphics World journal. In 1998, he joined OC Print Logic Technologies, as senior scientist. He worked there on various problem of image analysis dedicated to scanning and printing. After ten years of research work on image processing and computer graphics problems in several industrial companies, he joined the Informatics Department of ESIEE, Paris in 2002, where he is a professor and a member of the Laboratoire d’Informatique Gaspard Monge, Université Paris-Est Marne-la-Vallée. His current research interest is discrete mathematical morphology and discrete optimization.

# Supplementary material for paper *Tree-Based Morse Regions: A Topological Approach to Local Feature Detection*

Yongchao Xu, Pascal Monasse, Thierry Géraud, and Laurent Najman

## I. INTRODUCTION

In this supplementary material document, we show some complementary experimental results for the paper *Tree-Based Morse Regions: A Topological Approach to Local Feature Detection*. They are listed as below:

- In Section II, we show the result of the repeatability test (including DoG) on “Graffiti” sequence.
- In Section III, we show the comparison of features (along with the scales) extracted by different methods on a low contrast image and its contrast-enhanced one. This pair of images are shown in Fig. 6 of the paper.
- In Section IV, we show the location of detected points of images of CD-covers used in Fig. 11 and Fig. 13. These points are used in image registration experiments. The pairs of matched points used to estimate the homography are also illustrated.
- In Section V, we show some supplementary results of quantitative benchmark of registration experiments on Mikolajczyk dataset.

## II. REPEATABILITY TEST ON “GRAFFITI” SEQUENCE

We include the repeatability test on the “Graffiti” sequence in Mikolajczyk’s dataset. First of all, the number of detected points on each image of the sequence is given in Table I. TBMR extracts more points than MSER, but slightly less than Harris-Affine and Hessian-Affine, and DoG detects many more points than the others. The result of repeatability test (including DoG for which the round disks are used instead of squares) is shown in Fig. 1. The same observation is obtained as for the “Wall” sequence (with also the viewpoint change) shown in Fig. 9 in the paper. For the “Graffiti” sequence, in general, TBMR achieves better repeatability score than Harris-Affine and Hessian-Affine. Compared to MSER, TBMR obtains many more correspondences with a small loss of repeatability score. Note that we use the Linux executable available in <http://www.robots.ox.ac.uk/~vgg/research/affine/detectors.html> for the MSER. The public implementation of MSER in VLFeat [36] <http://www.vlfeat.org/> and OpenCV achieves results worse than the ones of the binary. DoG detects many points with a good repeatability score for small change of viewpoint. When the viewpoint change is large, the repeatability score decreases significantly, and from 50 degrees of viewpoint change. DoG does not give any correspondence. This explains the better performance in the homography experiments shown in Fig. 10 in the paper.

Methods	img1	img2	img3	img4	img5	img6
TBMR	1200	1285	1384	1588	1670	1886
MSER	547	633	700	707	798	925
Harris-Affine	1753	2059	2215	2065	2195	1888
Hessian-Affine	2510	2848	2782	2451	2385	1898
DoG	3198	3586	3911	3952	4299	4920

TABLE I. Number of points detected by different methods for each image of the “Graffiti” sequence. TBMR extracts more points than MSER, less than Harris-Affine and Hessian-Affine. DoG extracts many more points than the other methods.

Methods	Original image	contrast-enhanced image
TBMR	2978	2347
MSER	1025	1019
Harris-Affine	3119	1874
Hessian-Affine	2204	1474
DoG	7042	6768

TABLE II. Number of points detected by different methods on the low contrast image and its contrast-enhanced one. TBMR extracts more points than MSER, and comparable number of points with Harris-Affine and Hessian-Affine. DoG extracts many more points than the other methods.

## III. COMPARISON OF EXTRACTED LOCAL FEATURES ON A LOW CONTRAST IMAGE

We show the detected points with scales for the low contrast image and its contrast-enhanced one (as shown in Fig. 6 (a) of the paper). This original low contrast image and a contrast-enhanced one with significant increasing change of contrast are also shown in Fig. 2 (a). The number of extracted points by different methods are shown in Table II. The extracted points with scales of different methods are respectively shown in Fig. 2 (b), Fig. 3, and Fig. 4. As shown in Fig. 2, the points in mostly uniform regions shown in Fig. 5 (f) of the paper are actually points with a different large scale (the ellipses). MSER, DoG, Harris-Affine, and Hessian-Affine detect very few points in the area of low contrast (e.g., the body of the deer sculpture). By increasing the contrast, they detect some points. TBMR is perfectly insensitive to contrast change, up to quantization effects.

## IV. LOCAL FEATURE DETECTION FOR TWO IMAGES OF CD-COVERS

We show in Fig. 5 and Fig. 6 the locations of the interest points extracted by TBMR and DoG for the two CD-covers images: 007 and 010 used respectively in Fig. 11 and Fig. 13 of the paper. The number of interest points detected by these two methods is shown in Table III. For the reference images, TBMR detects fewer points than DoG. For the target images



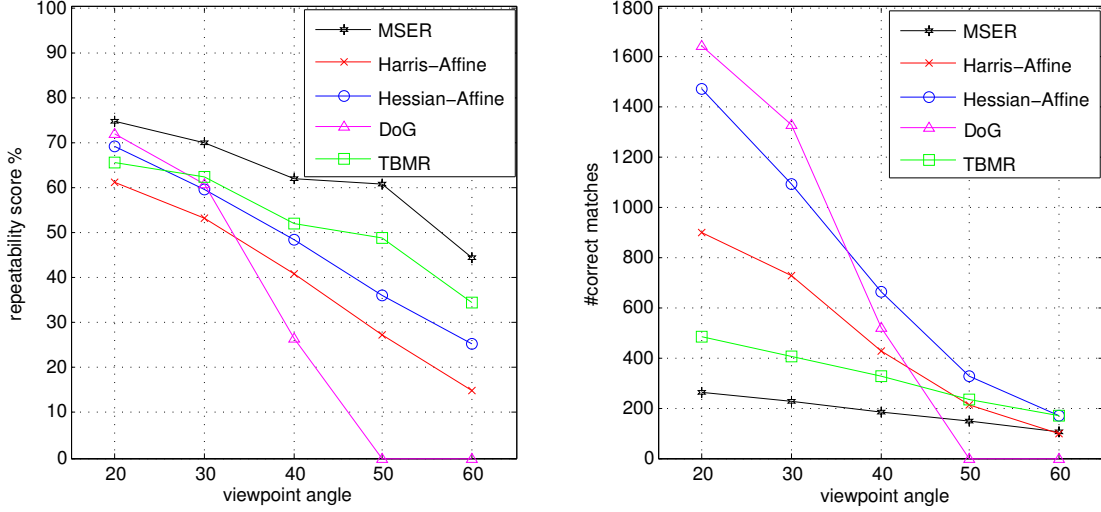


Fig. 1. Repeatability score (left) and number of correspondences (right) for the “Graffiti” sequence. TBMR is more robust with respect to viewpoint change, especially for strong viewpoint change.

Methods	Reference 007	Canon 007	Reference 010	Palm 010
TBMR	262	746	475	1581
DoG	1038	132	1017	537

TABLE III. Number of interest points detected by TBMR and DoG on two pairs of CD-covers images used respectively in Fig. 11 and Fig. 13.

taken with Canon and Palm cameras, TBMR extracts more points, but many of them lie in the area outside the cd covers. When we compute the homography between the two pairs of images, only the points extracted in reference images (inside the cd covers) are used. The registration based on TBMR makes use of fewer points than DoG. For the pair of images 007 of the CD-covers, both TBMR and DoG gives a registration result. The inlier matched pairs of points are given in Fig. 7. In fact, all the matched pairs of points given by TBMR are used as inliers to estimate the homography. One pair of matched points given by DoG is considered as outlier, which is not shown here. For the pair of images 010 of the CD-covers, all the matched pairs of points given by TBMR (shown in Fig. 8) are used as inliers to estimate the homography. DoG fails to estimate the homography, all the matched pairs of points are considered as outliers (see Fig. 8). Actually, four of the six DoG correspondences look correct. Although this is exactly the bare minimum required for a homography, there is no correct other correspondence to check its adequacy.

## V. IMAGE REGISTRATION ON MIKOLAJCZYK DATASET

We show some extra distributions of distances in image registration by homography on the dataset of Mikolajczyk *et al.* [6]. We compare mainly to MSER and DoG.

The results for the “Wall” sequence with viewpoint changes are shown in Fig. 9. Harris-Affine, Hessian-Affine, and DoG fail registering pairs (1,6); TBMR performs similar with MSER on all the pairs and similar with Harris-Affine, and Hessian on all other pairs than (1,6). In general, TBMR performs better than DoG on all the other pairs.

The results for the “Bark” sequence with scale changes are shown in Fig. 10. TBMR performs better than MSER for all

the pairs. Except for the pair (1,3) where TBMR is on par with DoG, TBMR performs better than DoG. Harris-Affine fails registering pair (1,6). In general, TBMR performs better than Hessian-Affine.

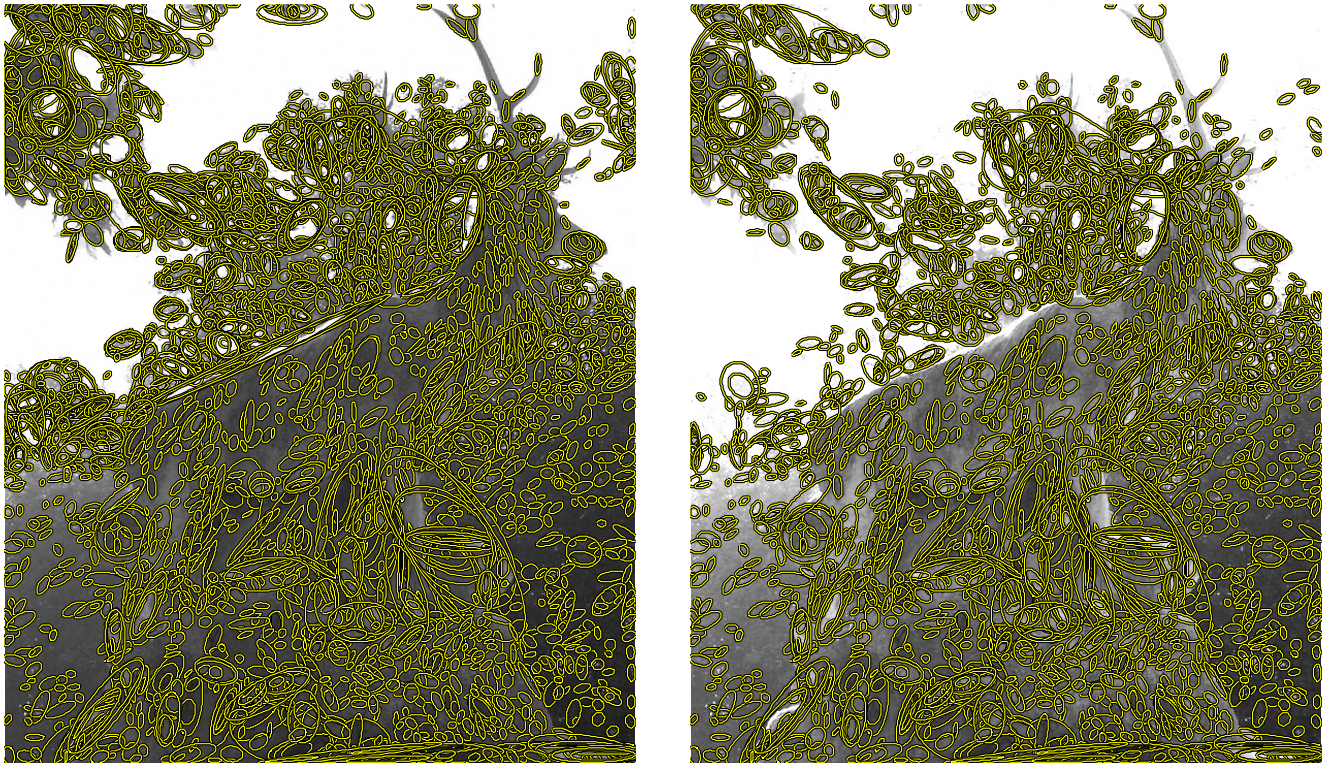
The results for the “Trees” sequence with different blur effects are shown in Fig. 11. TBMR performs better than MSER for all the pairs. Besides, MSER fails registering the pair (1,6). In general, TBMR also performs better than Harris-Affine and Hessian-Affine. When the blur is not very important (the pairs 1-2, 1-3, and 1-4), TBMR achieves similar results with DoG. But when the blur is very important, DoG performs better. This also explains the better performance of DoG in the registration experiments for “Palm” in Table I. However, as shown in [45], the multi-resolution detection improves the performance of MSER under blur. We would expect the same improvements by applying a multi-resolution analysis.

The results for the “Leuven” sequence with different blur amounts are shown in Fig. 12. In general, TBMR performs better than MSER, Harris-Affine, and Hessian-Affine. However, DoG performs better than TBMR: As all pairs of images cover almost the same scene and there are always some parts of the image having good contrast, even for the last image (img6) of the sequence, the change of contrast does not have a strong impact on the homography estimation for DoG.





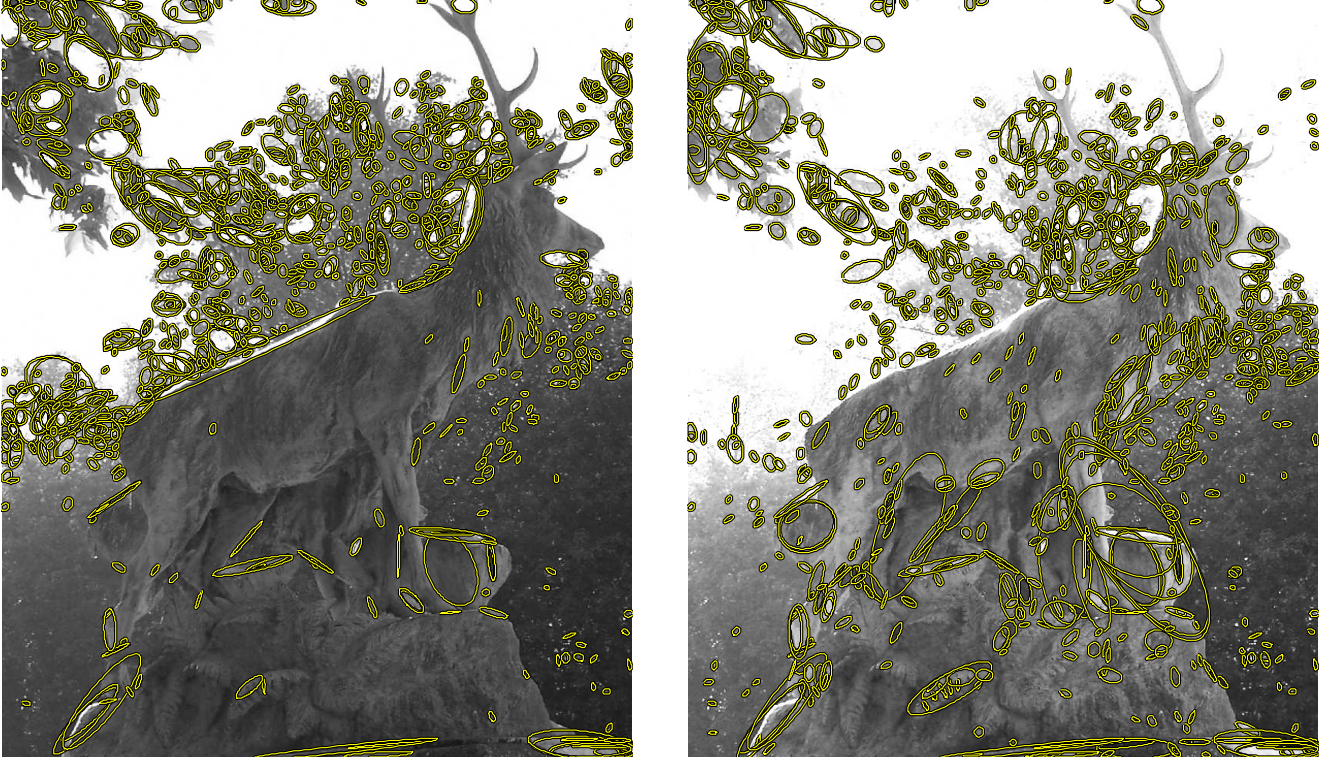
(a) Original low contrast image (left) and contrast-enhanced image (right).



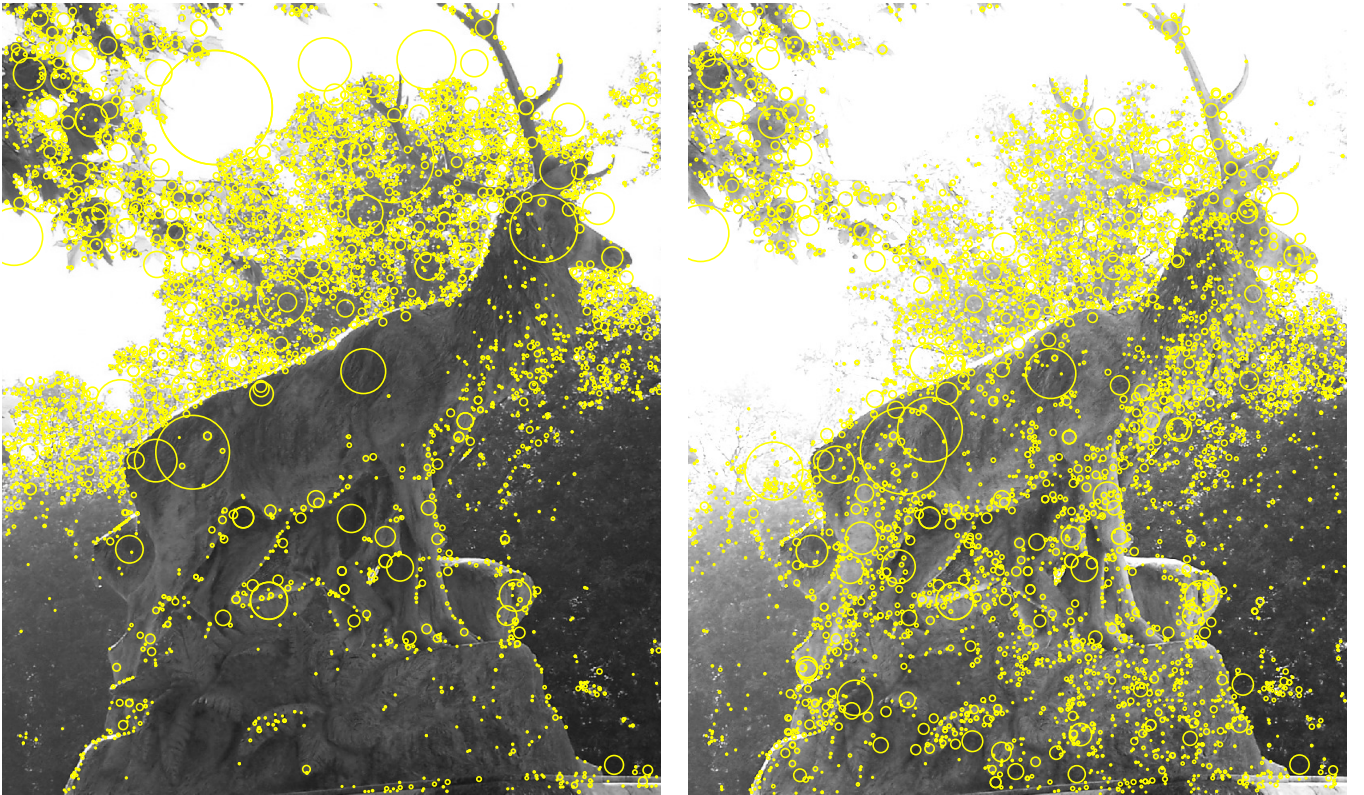
(b) Local features extracted by TBMR on original images (left) and contrast-enhanced image (right).

Fig. 2. Local features extracted by TBMR on low contrast image and contrast-enhanced image. The points in mostly uniform regions shown in Fig. 5 (f) of the paper are actually points with a different large scale (the ellipses). TBMRs are perfectly insensitive to contrast change, up to quantization effects.





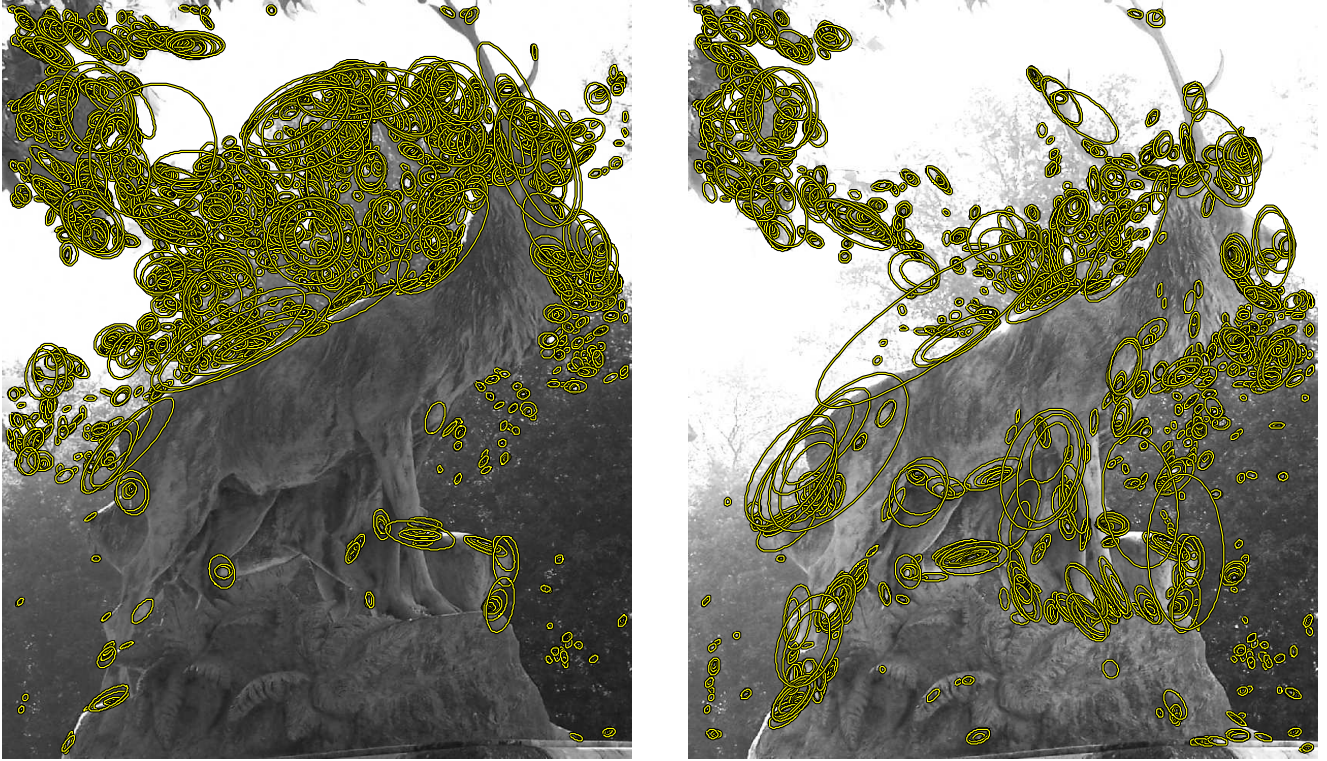
(a) Local features extracted by MSER on original images (left) and contrast-enhanced image (right).



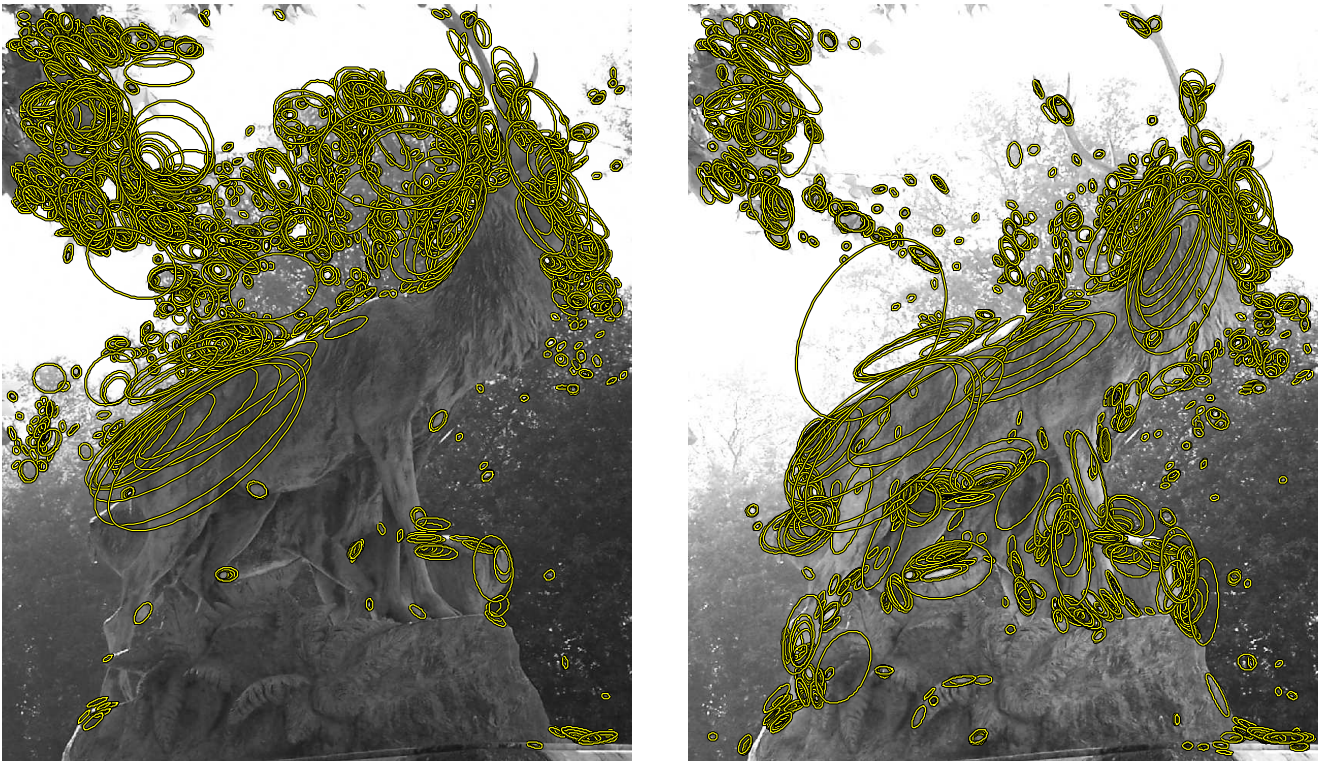
(b) Local features extracted by DoG on original images (left) and contrast-enhanced image (right).

Fig. 3. Local features extracted by MSER (top) and DoG (Down) on low contrast image and contrast-enhanced image. To better visualize the patches extracted by DoG, the round disks with diameter being the size of square patches divided by 6 are shown. MSER and DoG detect very few points in the area of low contrast (e.g., the body of the deer sculpture). By increasing the contrast, MSER and DoG detect some points.





(a) Local features extracted by Harris-Affine on original images (left) and contrast-enhanced image (right).



(b) Local features extracted by Hessian-Affine on original images (left) and contrast-enhanced image (right).

Fig. 4. Local features extracted by Harris-Affine (top) and Hessian-Affine (Down) on low contrast image and contrast-enhanced image. They both detect very few points in the area of low contrast (e.g., the body of the deer sculpture). By increasing the contrast, they detect some points.



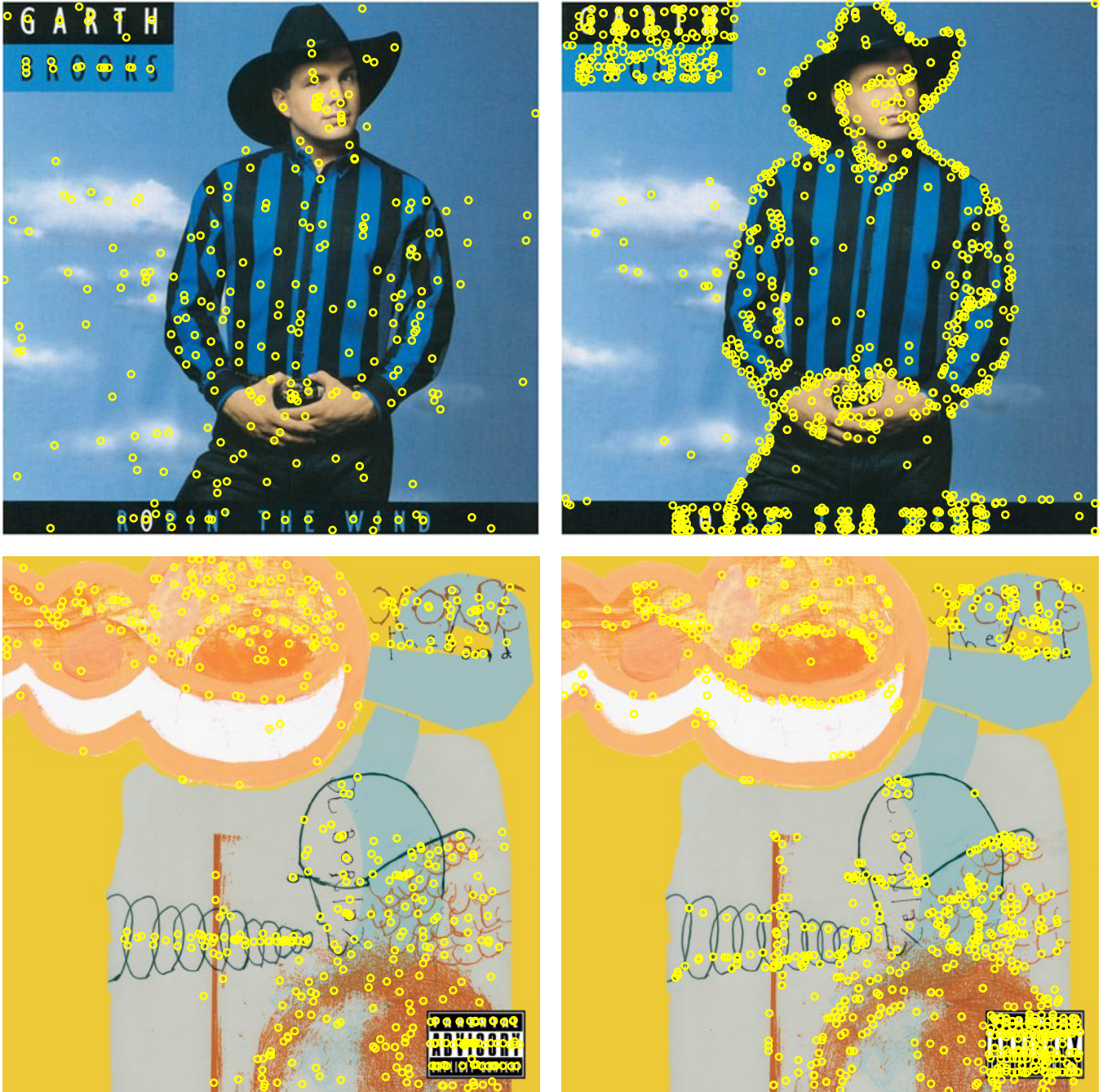


Fig. 5. Locations of extracted interest points detected by TBMR (left) and DoG (right) on the two reference images of CD-covers. TBMR detects fewer points than DoG.





Fig. 6. Locations of extracted interest points detected by TBMR (left) and DoG (right) on the two target images taken by respectively Canon (top) and Palm (down) cameras of CD-covers. TBMR detects more points than DoG. Many of them lie in the area outside of the cd covers. In the area of low contrast, TBMR extracts more points than DoG.



Fig. 7. Inlier matched pairs of points using respectively TBMR (top) and DoG (down). TBMR has more inlier matched pairs of points than DoG. In fact, all the matched pairs of points using TBMR are considered as the inliers to estimate the homography. One pair of matched points using DoG is considered as outliers in the homography estimation.



Fig. 8. Matched pairs of points using respectively TBMR (top) and DoG (down). All the matched pairs of points using TBMR are used as inliers to estimate the homography. Using DoG fails to estimate the homography: all the matched pairs of points are considered as outliers. Though four pairs are correct, there is no other to check automatically the resulting homography.



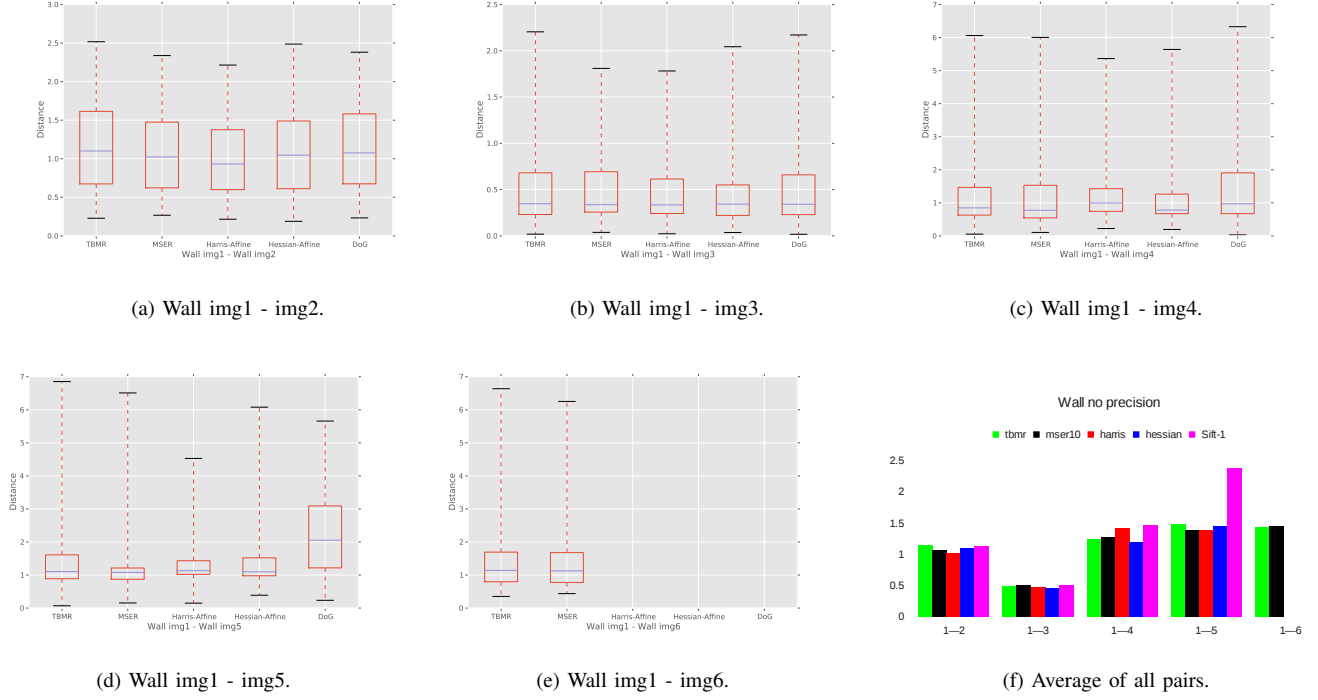


Fig. 9. Distributions of errors in image registration by homography on “Wall” image pairs, which comes with a ground truth. Harris-Affine, Hessian-Affine, and DoG fails registering pairs (1,6); TBMR performs similar with MSER on all the pairs and similar with Harris-Affine, and Hessian on all other pairs than (1,6). In general, TBMR performs better than DoG on all the other pairs.

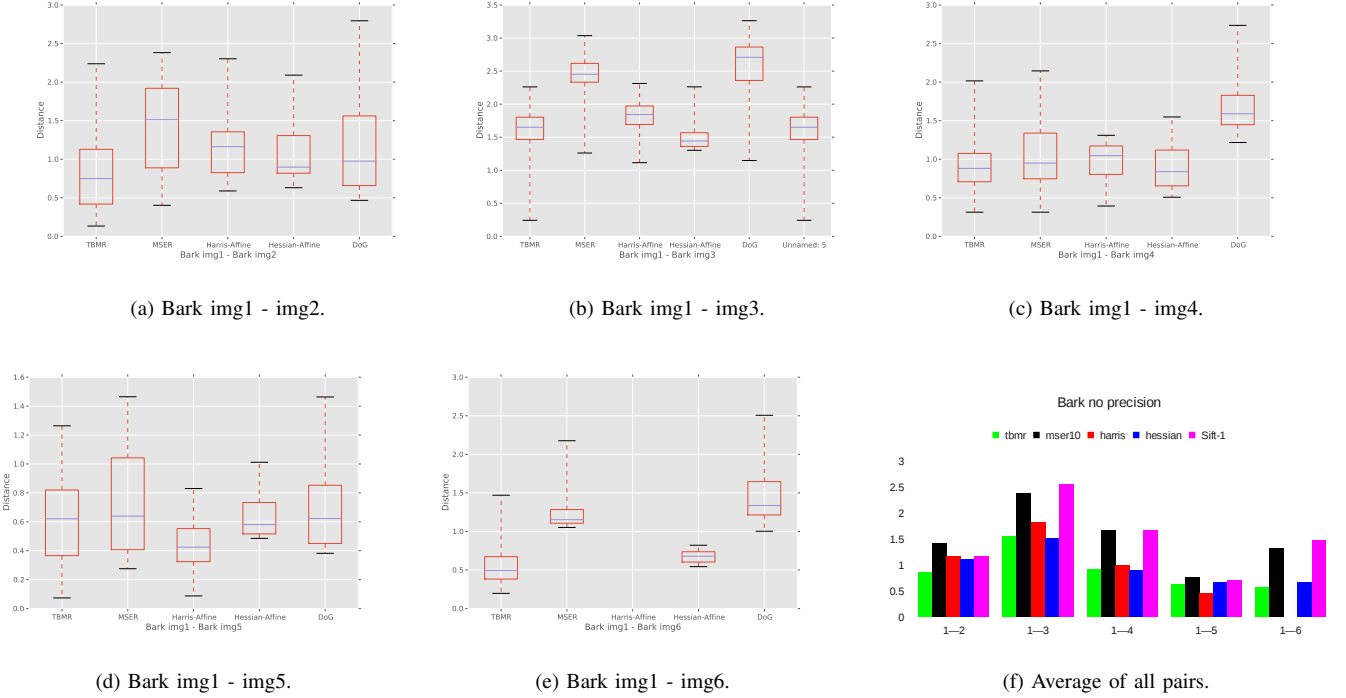


Fig. 10. Distributions of errors in image registration by homography on “Bark” image pairs, which comes with a ground truth. TBMR performs better than MSER for all the pairs. Except for the pair (1,3) where TBMR similar with DoG, TBMR performs better than DoG. Harris-Affine fails registering pair (1,6). And in general, TBMR performs better than Hessian-Affine.

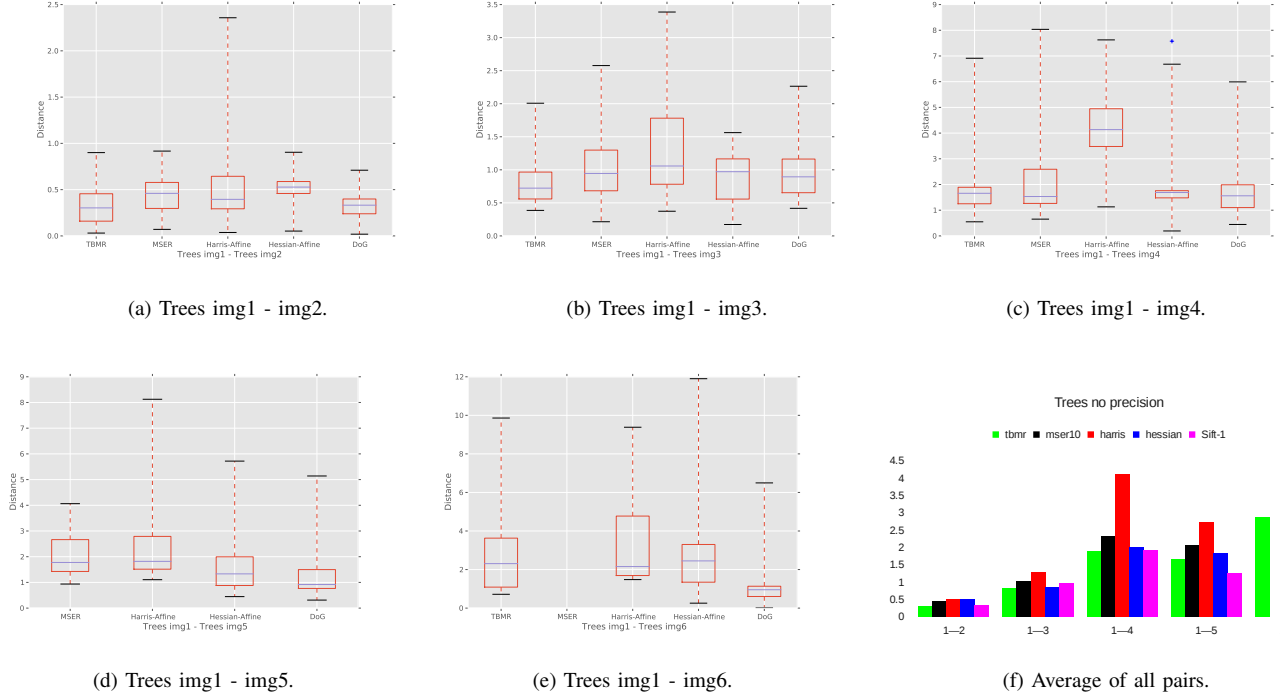


Fig. 11. Distributions of errors in image registration by homography on “Trees” image pairs, which comes with a ground truth. TBMR performs better than MSER for all the pairs. Besides, MSER fails registering the pair (1,6). In general, TBMR also performs better than Harris-Affine and Hessian-Affine. When the blur is not very important (the pairs 1-2, 1-3, and 1-4), TBMR achieves similar results with DoG. But when the blur is very important, DoG performs better.

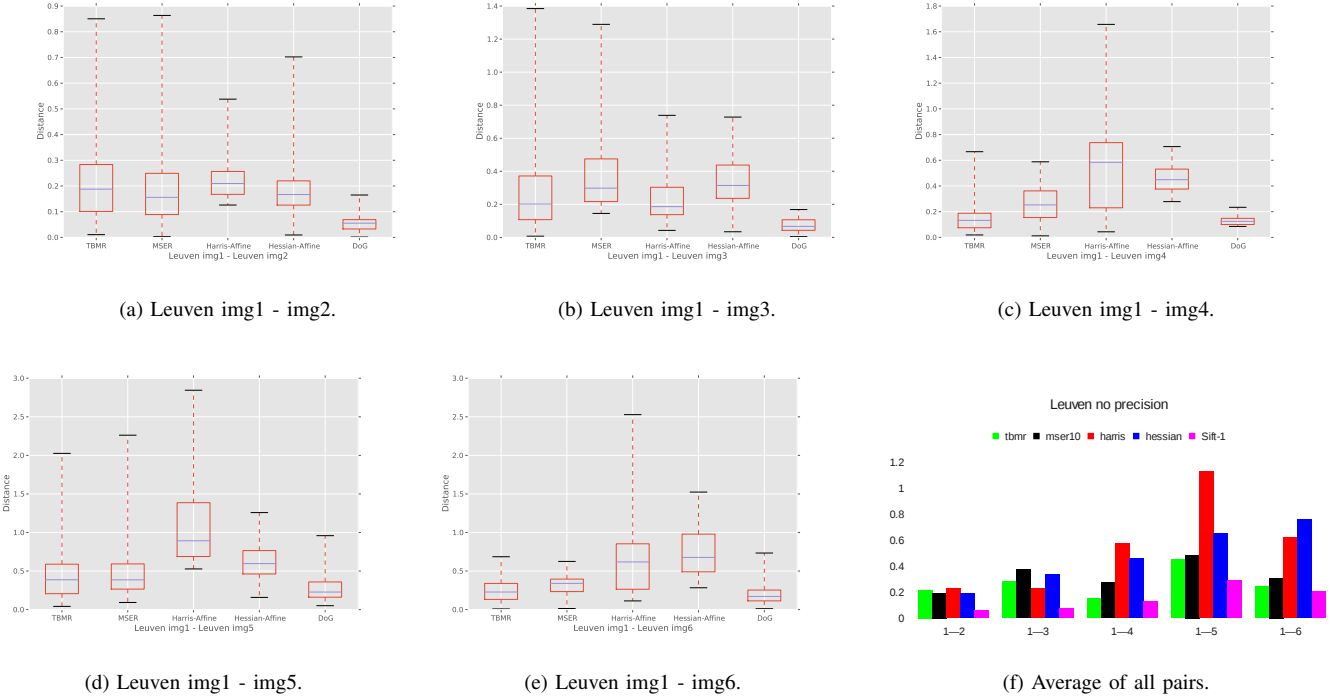


Fig. 12. Distributions of errors in image registration by homography on “Leuven” image pairs, which comes with a ground truth. In general, TBMR performs better than MSER, Harris-Affine, and Hessian-Affine. However, DoG performs better than TBMR: As all pairs of images cover almost the same scene and there are always some parts of the image having good contrast, even for the last image (img6), the change of contrast does not have a strong impact on the homography estimation for DoG.